

# Quantifying Synergistic Information

Thesis by  
Virgil Griffith

In Partial Fulfillment of the Requirements  
for the Degree of  
Doctor of Philosophy



California Institute of Technology  
Pasadena, California

2013  
(Submitted January 31, 2014)

© 2013

Virgil Griffith

All Rights Reserved

This thesis would not exist were it not for many people. Some of the prominent ones are Douglas R. Hofstadter, for the inspiration; Giulio Tononi, for the theory; Christof Koch, for the ambition to explore new waters; and the DOE CSGF, for the funding to work in peace.

# Acknowledgments

I wish to thank Christof Koch, Suzannah A. Fraker, Paul Williams, Mark Burgin, Tracey Ho, Edwin K. P. Chong, Christopher J. Ellison, Ryan G. James, Jim Beck, Shuki Bruck, and Pietro Perona.



# Abstract

Within the microcosm of information theory, I explore what it means for a system to be functionally irreducible. This is operationalized as quantifying the extent to which cooperative or “synergistic” effects enable random variables  $X_1, \dots, X_n$  to predict (have mutual information about) a single target random variable  $Y$ . In Chapter 1, we introduce the problem with some emblematic examples. In Chapter 2, we show how six different measures from the existing literature fail to quantify this notion of synergistic mutual information. In Chapter 3, we take a step towards a measure of synergy which yields the first nontrivial lowerbound on synergistic mutual information. In Chapter 4, we find that synergy is but the weakest notion of a broader concept of *irreducibility*. In Chapter 5, we apply our results from Chapters 3 and 4 towards grounding Giulio Tononi’s ambitious  $\phi$  measure, which attempts to quantify the magnitude of consciousness experience.

# Contents

<b>Acknowledgments</b>	<b>iv</b>
<b>Abstract</b>	<b>v</b>
<b>I Introducing the Problem</b>	<b>4</b>
<b>1 What is Synergy?</b>	<b>5</b>
1.1 Notation and PI-diagrams . . . . .	5
1.1.1 Understanding PI-diagrams . . . . .	6
1.2 Information can be redundant, unique, or synergistic . . . . .	7
1.2.1 Example Rdn: Redundant information . . . . .	8
1.2.2 Example Unq: Unique information . . . . .	8
1.2.3 Example Xor: Synergistic information . . . . .	8
<b>2 Six Prior Measures of Synergy</b>	<b>10</b>
2.1 Definitions . . . . .	10
2.1.1 Multivariate Mutual Information: $\text{MMI}(X_1; \dots; X_n; Y)$ . . . . .	10
2.1.2 Interaction Information: $\mathcal{II}(X_1; \dots; X_n; Y)$ . . . . .	10
2.1.3 WholeMinusSum synergy: $\text{WMS}(\mathbf{X}; Y)$ . . . . .	11
2.1.4 WholeMinusPartitionSum: $\text{WMPS}(\mathbf{X}; Y)$ . . . . .	12
2.1.5 $I_{\max}$ synergy: $\mathcal{S}_{\max}(\mathbf{X}; Y)$ . . . . .	13
2.1.6 Correlational importance: $\Delta I(\mathbf{X}; Y)$ . . . . .	14
2.2 The six prior measures are not equivalent . . . . .	15
2.3 Counter-intuitive behaviors of the six prior measures . . . . .	15
2.3.1 $I_{\max}$ synergy: $\mathcal{S}_{\max}$ . . . . .	15
2.3.2 $\mathcal{S}_{\text{MMI}}$ , $\mathcal{II}$ , $\text{WMS}$ , $\text{WMPS}$ . . . . .	16
2.3.3 Correlational importance: $\Delta I$ . . . . .	17
2.A Algebraic simplification of $\Delta I$ . . . . .	20

<b>II</b>	<b>Making Progress</b>	<b>22</b>
<b>3</b>	<b>First Nontrivial Lowerbound on Synergy</b>	<b>23</b>
3.1	Introduction . . . . .	23
3.2	Two examples elucidating desired properties for synergy . . . . .	24
3.2.1	XorDuplicate: Synergy is invariant to duplicating a predictor . . . . .	24
3.2.2	XorLoses: Adding a new predictor can decrease synergy . . . . .	25
3.3	Preliminaries . . . . .	25
3.3.1	Informational Partial Order and Equivalence . . . . .	25
3.3.2	Information Lattice . . . . .	27
3.3.3	Invariance and Monotonicity of Entropy . . . . .	28
3.3.4	Desired Properties of Intersection Information . . . . .	28
3.4	Candidate Intersection Information for Zero-Error Information . . . . .	30
3.4.1	Zero-Error Information . . . . .	30
3.4.2	Intersection Information for Zero-Error Information . . . . .	31
3.5	Candidate Intersection Information for Shannon Information . . . . .	31
3.6	Three Examples Comparing $I_{\min}$ and $I_{\lambda}$ . . . . .	33
3.7	Negative synergy and state-dependent ( <b>GP</b> ) . . . . .	35
3.7.1	Consequences of state-dependent ( <b>GP</b> ) . . . . .	37
3.8	Conclusion and Path Forward . . . . .	37
3.A	Algorithm for Computing Common Random Variable . . . . .	40
3.B	Algorithm for Computing $I_{\lambda}$ . . . . .	40
3.C	Lemmas and Proofs . . . . .	40
3.C.1	Lemmas on Desired Properties . . . . .	40
3.C.2	Properties of $I_{\lambda}^0$ . . . . .	41
3.C.3	Properties of $I_{\lambda}$ . . . . .	44
3.D	Miscellaneous Results . . . . .	47
3.E	Misc Figures . . . . .	49
<b>4</b>	<b>Irreducibility is Minimum Synergy among Parts</b>	<b>50</b>
4.1	Introduction . . . . .	50
4.1.1	Notation . . . . .	50
4.2	Four common notions of irreducibility . . . . .	51
4.3	Quantifying the four notions of irreducibility . . . . .	52
4.3.1	Information beyond the Elements . . . . .	53
4.3.2	Information beyond Disjoint Parts: <b>lbDp</b> . . . . .	53
4.3.3	Information beyond Two Parts: <b>lb2p</b> . . . . .	53

4.3.4	Information beyond All Parts: lbAp . . . . .	54
4.4	Exemplary Binary Circuits . . . . .	54
4.4.1	XorUnique: Irreducible to elements, yet reducible to a partition . . . . .	55
4.4.2	DoubleXor: Irreducible to a partition, yet reducible to a pair . . . . .	55
4.4.3	TripleXor: Irreducible to a pair of components, yet still reducible . . . . .	56
4.4.4	Parity: Complete irreducibility . . . . .	57
4.5	Conclusion . . . . .	57
4.A	Joint distributions for DoubleXor and TripleXor . . . . .	62
4.B	Proofs . . . . .	62

### III Applications 66

5	Improving the $\phi$ Measure <span style="float: right;">67</span>
5.1	Introduction . . . . . 67
5.2	Preliminaries . . . . . 67
5.2.1	Notation . . . . . 67
5.2.2	Model assumptions . . . . . 68
5.3	How $\phi$ works . . . . . 68
5.3.1	Stateless $\phi$ is $\langle \phi \rangle$ . . . . . 70
5.4	Room for improvement in $\phi$ . . . . . 70
5.5	A Novel Measure of Irreducibility to a Partition . . . . . 74
5.5.1	Stateless $\psi$ is $\langle \psi \rangle$ . . . . . 75
5.6	Contrasting $\psi$ versus $\phi$ . . . . . 76
5.7	Conclusion . . . . . 77
5.A	Reading the network diagrams . . . . . 80
5.B	Necessary proofs . . . . . 83
5.B.1	Proof that the max union of bipartitions covers all partitions . . . . . 83
5.B.2	Bounds on $\psi(X_1, \dots, X_n : y)$ . . . . . 85
5.B.3	Bounds on $\langle \psi \rangle(X_1, \dots, X_n : Y)$ . . . . . 89
5.C	Definition of intrinsic $\text{ei}(y/\mathbf{P})$ a.k.a. “perturbing the wires” . . . . . 91
5.D	Misc proofs . . . . . 93
5.E	Setting $t = 1$ without loss of generality . . . . . 94

### Bibliography 95

## Part I

# Introducing the Problem

# Chapter 1

## What is Synergy?

The prior literature [24, 30, 1, 6, 19, 36] has termed several distinct concepts as “synergy”. We define synergy as a special case of *irreducibility*—specifically, synergy is irreducibility to atomic elements. By definition, a group of two or more agents synergistically perform a task if and only if the performance of that task decreases when the agents work “separately”, or in parallel isolation. It is important to remember that it is the collective *action* that is irreducible, not the agents themselves. A concrete example of irreducibility is the “agents” hydrogen and oxygen working to extinguish fire. Even when  $H_2$  and  $O_2$  are both present in the same container, if working separately neither extinguishes fire (on the contrary, fire grows!). But hydrogen and oxygen fused or “grouped” into a single entity,  $H_2O$ , readily extinguishes fire.

The concept of synergy spans many fields and theoretically could be applied to any non-subadditive function. But within the confines of Shannon information theory, synergy—or more formally, *synergistic information*—is a property of a set of  $n$  random variables  $\mathbf{X} = \{X_1, X_2, \dots, X_n\}$  cooperating to predict, that is, reduce the uncertainty of, a single target random variable  $Y$ .

### 1.1 Notation and PI-diagrams

We use the following notation throughout. Let

$n$ : The number of predictors  $X_1, X_2, \dots, X_n$ .  $n \geq 2$ .

$X_{1\dots n}$ : The joint random variable (cartesian product) of all  $n$  predictors  $X_1 X_2 \dots X_n$ .

$X_i$ : The  $i$ ’th predictor random variable (r.v.).  $1 \leq i \leq n$ .

$\mathbf{X}$ : The *set* of all  $n$  predictors  $\{X_1, X_2, \dots, X_n\}$ .

$Y$ : The *target r.v.* to be predicted.

$y$ : A particular state of the target r.v.  $Y$ .

All random variables are discrete, all logarithms are  $\log_2$ , and all calculations are in *bits*. Entropy and mutual information are as defined by [9],  $H(X) \equiv \sum_{x \in X} \Pr(x) \log \frac{1}{\Pr(x)}$ , as well as  $I(X:Y) \equiv \sum_{x,y} \Pr(x,y) \log \frac{\Pr(x,y)}{\Pr(x)\Pr(y)}$ .

### 1.1.1 Understanding PI-diagrams

Partial information diagrams (PI-diagrams), introduced by [36], extend Venn diagrams to properly represent synergy. Their framework has been invaluable to the evolution of our thinking on synergy.

A PI-diagram is composed of nonnegative *partial information regions* (PI-regions). Unlike the standard Venn entropy diagram in which the sum of all regions is the joint entropy  $H(X_{1\dots n}, Y)$ , in PI-diagrams the sum of all regions (i.e. the space of the PI-diagram) is the mutual information  $I(X_{1\dots n}:Y)$ . PI-diagrams are immensely helpful in understanding how the mutual information  $I(X_{1\dots n}:Y)$  is distributed across the coalitions and singletons of  $\mathbf{X}$ .<sup>1</sup>

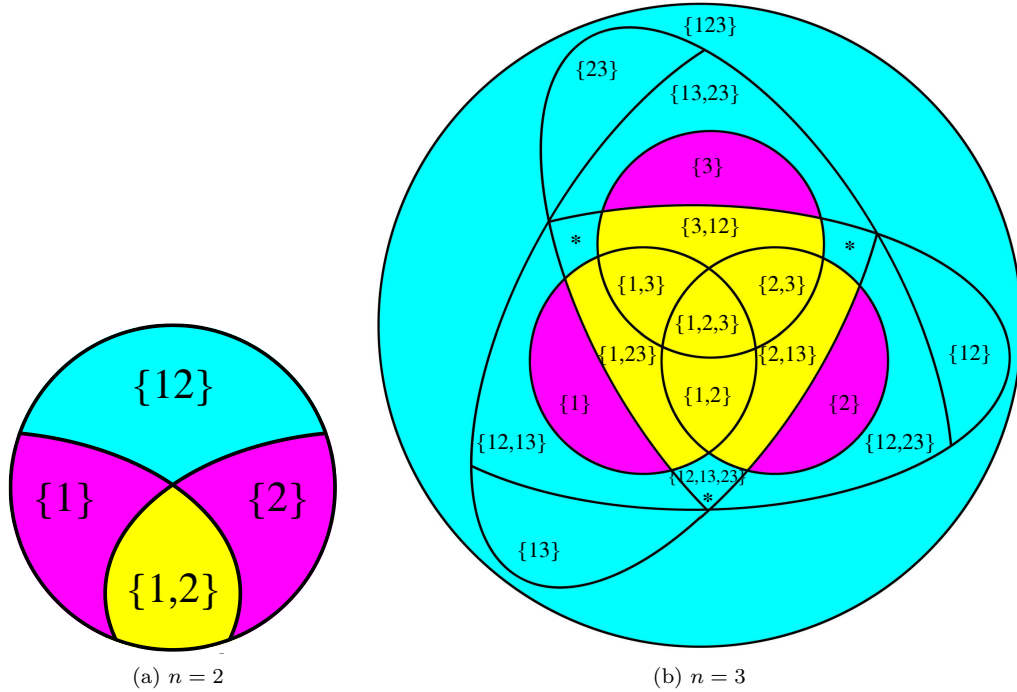


Figure 1.1: PI-diagrams for two and three predictors. Each PI-region represents nonnegative information about  $Y$ . A PI-region’s color represents whether its information is redundant (yellow), unique (magenta), or synergistic (cyan). To preserve symmetry, the PI-region “ $\{12, 13, 23\}$ ” is displayed as three separate regions each marked with a “\*”. All three \*-regions should be treated as though they are a single region.

**How to read PI-diagrams.** Each PI-region is uniquely identified by its “set notation” where each element is denoted solely by the predictors’ indices. For example, in the PI-diagram for  $n = 2$

<sup>1</sup>Formally, how the mutual information is distributed across the set of all nonempty antichains on the powerset of  $\mathbf{X}$ [35].

(Figure 1.1a):  $\{1\}$  is the information about  $Y$  only  $X_1$  carries (likewise  $\{2\}$  is the information only  $X_2$  carries);  $\{1,2\}$  is the information about  $Y$  that  $X_1$  as well as  $X_2$  carries, while  $\{12\}$  is the information about  $Y$  that is specified only by the coalition (joint random variable)  $X_1X_2$ . For getting used to this way of thinking, common informational quantities are represented by colored regions in Figure 5.5.

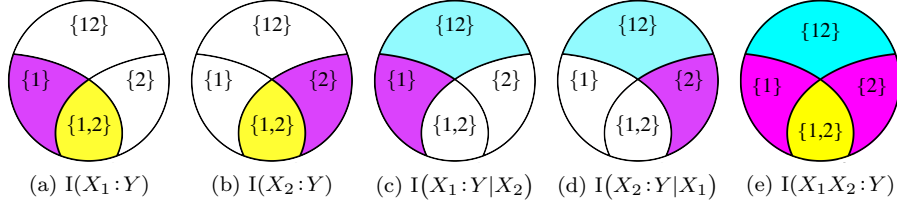


Figure 1.2: PI-diagrams for  $n = 2$  representing standard informational quantities.

The general structure of a PI-diagram becomes clearer after examining the PI-diagram for  $n = 3$  (Figure 1.1b). All PI-regions from  $n = 2$  are again present. Each predictor ( $X_1, X_2, X_3$ ) can carry unique information (regions labeled  $\{1\}, \{2\}, \{3\}$ ), carry information redundantly with another predictor ( $\{1,2\}, \{1,3\}, \{2,3\}$ ), or specify information through a coalition with another predictor ( $\{12\}, \{13\}, \{23\}$ ). New in  $n = 3$  is information carried by all three predictors ( $\{1,2,3\}$ ) as well as information specified through a three-way coalition ( $\{123\}$ ). Intriguingly, for three predictors, information can be provided by a coalition as well as a singleton ( $\{1,23\}, \{2,13\}, \{3,12\}$ ) or specified by multiple coalitions ( $\{12,13\}, \{12,23\}, \{13,23\}, \{12,13,23\}$ ).

## 1.2 Information can be redundant, unique, or synergistic

Each PI-region represents an irreducible nonnegative slice of the mutual information  $I(X_{1\dots n}:Y)$  that is either:

1. **Redundant.** Information carried by a singleton predictor as well as available somewhere else.  
For  $n = 2$ :  $\{1,2\}$ . For  $n = 3$ :  $\{1,2\}, \{1,3\}, \{2,3\}, \{1,2,3\}, \{1,23\}, \{2,13\}, \{3,12\}$ .
2. **Unique.** Information carried by exactly one singleton predictor and available nowhere else.  
For  $n = 2$ :  $\{1\}, \{2\}$ . For  $n = 3$ :  $\{1\}, \{2\}, \{3\}$ .
3. **Synergistic.** Any and all information in  $I(X_{1\dots n}:Y)$  that is not carried by a singleton predictor.  
 $n = 2$ :  $\{12\}$ . For  $n = 3$ :  $\{12\}, \{13\}, \{23\}, \{123\}, \{12,13\}, \{12,23\}, \{13,23\}, \{12,13,23\}$ .

Although a single PI-region is either redundant, unique, or synergistic, a single state of the target can have any combination of positive PI-regions, i.e. a single state of the target can convey redundant, unique, and synergistic information. This surprising fact is demonstrated in Figure 3.4.



### 1.2.1 Example Rdn: Redundant information

If  $X_1$  and  $X_2$  carry some of the same information<sup>2</sup> (reduce the same uncertainty) about  $Y$ , then we say the set  $\mathbf{X} = \{X_1, X_2\}$  has some *redundant information* about  $Y$ . Figure 1.3 illustrates a simple case of redundant information.  $Y$  has two equiprobable states:  $\mathbf{r}$  and  $\mathbf{R}$  ( $\mathbf{r}/\mathbf{R}$  for “redundant bit”). Examining  $X_1$  or  $X_2$  identically specifies one bit of  $Y$ , thus we say set  $\mathbf{X} = \{X_1, X_2\}$  has one bit of redundant information about  $Y$ .

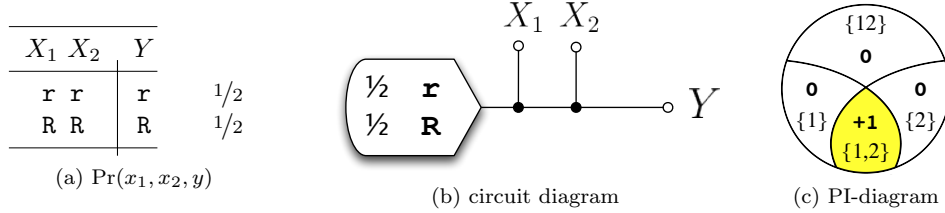


Figure 1.3: Example RDN. Figure 1.3a shows the joint distribution of r.v.’s  $X_1$ ,  $X_2$ , and  $Y$ , and the joint probability  $\Pr(x_1, x_2, y)$  is along the right-hand side of (a), revealing that all three terms are fully correlated. Figure 1.3b represents the joint distribution as an electrical circuit. Figure 1.3c is the PI-diagram indicating that set  $\{X_1, X_2\}$  has 1 bit of redundant information about  $Y$ .  $I(X_1 X_2 : Y) = I(X_1 : Y) = I(X_2 : Y) = H(Y) = 1$  bit.

### 1.2.2 Example Unq: Unique information

Predictor  $X_i$  carries *unique information* about  $Y$  if and only if  $X_i$  specifies information about  $Y$  that is not specified by anything else (a singleton or coalition of the other  $n - 1$  predictors). Figure 1.4 illustrates a simple case of unique information.  $Y$  has four equiprobable states:  $\mathbf{ab}$ ,  $\mathbf{aB}$ ,  $\mathbf{Ab}$ , and  $\mathbf{AB}$ .  $X_1$  uniquely specifies bit  $\mathbf{a}/\mathbf{A}$ , and  $X_2$  uniquely specifies bit  $\mathbf{b}/\mathbf{B}$ . If we had instead labeled the  $Y$ -states: 0, 1, 2, and 3,  $X_1$  and  $X_2$  would still have strictly unique information about  $Y$ . The state of  $X_1$  would specify between  $\{0, 1\}$  and  $\{2, 3\}$ , and the state of  $X_2$  would specify between  $\{0, 2\}$  and  $\{1, 3\}$ —together fully specifying the state of  $Y$ .

### 1.2.3 Example Xor: Synergistic information

A set of predictors  $\mathbf{X} = \{X_1, \dots, X_n\}$  has synergistic information about  $Y$  if and only if the whole  $(X_1 \dots X_n)$  specifies information about  $Y$  that is not specified by any singleton predictor. The canonical example of synergistic information is the XOR-gate (Figure 1.5). In this example, the whole  $X_1 X_2$  fully specifies  $Y$ ,

$$I(X_1 X_2 : Y) = H(Y) = 1 \text{ bit},$$

<sup>2</sup> $X_1$  and  $X_2$  providing identical information about  $Y$  is different from providing the same *magnitude* of information about  $Y$ , i.e.  $I(X_1 : Y) = I(X_2 : Y)$ . Example UNQ (Figure 1.4) is an example where  $I(X_1 : Y) = I(X_2 : Y) = 1$  bit yet  $X_1$  and  $X_2$  specify “different bits” of  $Y$ . Providing the same magnitude of information about  $Y$  is neither necessary or sufficient for providing some identical information about  $Y$ .

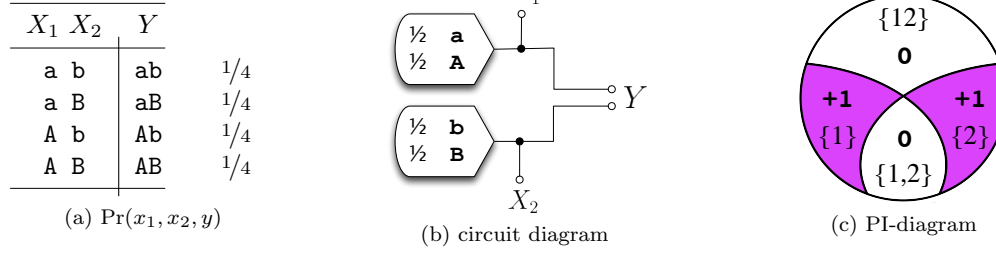


Figure 1.4: Example UNQ.  $X_1$  and  $X_2$  each uniquely specify a single bit of  $Y$ .  $I(X_1 X_2 : Y) = H(Y) = 2$  bits. The joint probability  $\Pr(x_1, x_2, y)$  is along the right-hand side of (a).

but the singletons  $X_1$  and  $X_2$  specify *nothing* about  $Y$ ,

$$I(X_1 : Y) = I(X_2 : Y) = 0 \text{ bits.}$$

With both  $X_1$  and  $X_2$  themselves having zero information about  $Y$ , we know that there can not be any redundant or unique information about  $Y$ , that the three PI-regions  $\{1\} = \{2\} = \{1, 2\} = 0$  bits. As the information between  $X_1 X_2$  and  $Y$  must come from somewhere, by elimination we conclude that  $X_1$  and  $X_2$  synergistically specify  $Y$ .

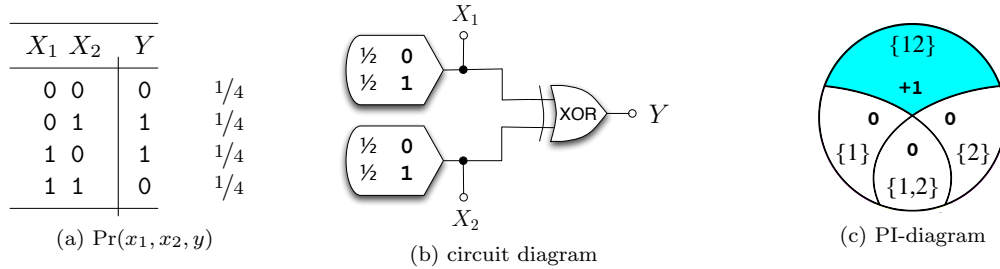


Figure 1.5: Example XOR.  $X_1$  and  $X_2$  synergistically specify  $Y$ .  $I(X_1 X_2 : Y) = H(Y) = 1$  bit. The joint probability  $\Pr(x_1, x_2, y)$  is along the right-hand side of (a).

## Chapter 2

# Six Prior Measures of Synergy

### 2.1 Definitions

#### 2.1.1 Multivariate Mutual Information: $\text{MMI}(X_1; \dots; X_n; Y)$

The first information-theoretic measure of synergy dates to 1954 from [24]. Inspired by Venn entropy diagrams, they defined the *multivariate mutual information* (MMI),  $\text{MMI}(X_1; \dots; X_n; Y) \equiv \sum_{T \subseteq \{X_1, \dots, X_n, Y\}} (-1)^{|T|+1} H(T)$ . Negative MMI was understood to be synergy. Therefore the MMI measure of synergy is,

$$\begin{aligned} \mathcal{S}_{\text{MMI}}(X_1; \dots; X_n; Y) &\equiv - \sum_{\mathbf{T} \subseteq \{X_1, \dots, X_n, Y\}} (-1)^{|\mathbf{T}|+1} H(T) \\ &= \sum_{\mathbf{T} \subseteq \{X_1, \dots, X_n, Y\}} (-1)^{|\mathbf{T}|} H(T) . \end{aligned} \tag{2.1}$$

#### 2.1.2 Interaction Information: $\mathcal{II}(X_1; \dots; X_n; Y)$

Interaction information ( $\mathcal{II}$ ), sometimes called the co-information, was introduced in [6] and tweaks MMI synergy measure. Although intended to measure informational “groupness” [6], Interaction Information is commonly interpreted as the magnitude of “information bound up in a set of variables, beyond that which is present in any subset of those variables.”<sup>1</sup>

Interaction Information among the  $n$  predictors and  $Y$  is defined as,

$$\begin{aligned} \mathcal{II}(X_1; \dots; X_n; Y) &\equiv (-1)^n \mathcal{S}_{\text{MMI}}(X_1; \dots; X_n; Y) \\ &= \sum_{\mathbf{T} \subseteq \{X_1, \dots, X_n, Y\}} (-1)^{n-|\mathbf{T}|} H(T) . \end{aligned} \tag{2.2}$$

---

<sup>1</sup>From [http://en.wikipedia.org/wiki/Interaction\\_information](http://en.wikipedia.org/wiki/Interaction_information).

Interaction Information is a signed measure where a positive value signifies synergy and a negative value signifies redundancy. Representing Interaction Information as a PI-diagram (Figure 2.1) reveals an intimidating imbroglio of added and subtracted PI-regions.

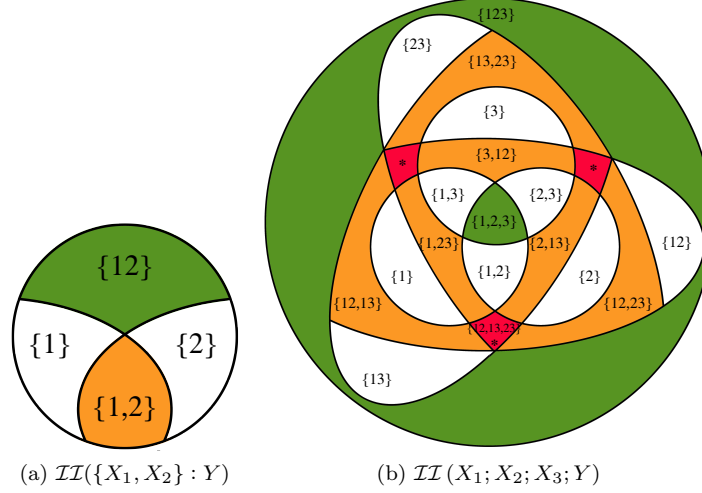


Figure 2.1: PI-diagrams illustrating interaction information for  $n = 2$  (left) and  $n = 3$  (right). The colors denote the added and subtracted PI-regions.  $\text{WMS}(\mathbf{X} : Y)$  is the green PI-region(s), minus the orange PI-region(s), minus two times any red PI-region.

### 2.1.3 WholeMinusSum synergy: $\text{WMS}(\mathbf{X} : Y)$

The earliest known sightings of bivariate WholeMinusSum synergy (WMS) are in [13, 12], with the general case in [11]. WholeMinusSum synergy is a signed measure where a positive value signifies synergy and a negative value signifies redundancy. WholeMinusSum synergy is defined by eq. (2.3) and interestingly reduces to eq. (2.5)—the difference of two *total correlations*.<sup>2</sup>

$$\text{WMS}(\mathbf{X} : Y) \equiv I(X_{1\dots n} : Y) - \sum_{i=1}^n I(X_i : Y) \quad (2.3)$$

$$= \sum_{i=1}^n H(X_i | Y) - H(X_{1\dots n} | Y) - \left[ \sum_{i=1}^n H(X_i) - H(X_{1\dots n}) \right] \quad (2.4)$$

$$= \text{TC}(X_1; \dots; X_n | Y) - \text{TC}(X_1; \dots; X_n) \quad (2.5)$$

Representing eq. (2.3) for  $n = 2$  as a PI-diagram (Figure 2.2a) reveals that WMS is the synergy between  $X_1$  and  $X_2$  *minus* their redundancy. Thus, when there is an equal magnitude of synergy and redundancy between  $X_1$  and  $X_2$ , WholeMinusSum synergy is *zero*—leading one to *erroneously*

<sup>2</sup> $\text{TC}(X_1; \dots; X_n) = -H(X_{1\dots n}) + \sum_{i=1}^n H(X_i)$  per [17].

conclude there is no synergy or redundancy present.<sup>3</sup>

The PI-diagram for  $n = 3$  (Figure 2.2b) reveals that WholeMinusSum double-subtracts PI-regions  $\{1,2\}$ ,  $\{1,3\}$ ,  $\{2,3\}$  and triple-subtracts PI-region  $\{1,2,3\}$ , revealing that for  $n > 2$   $\text{WMS}(\mathbf{X} : Y)$  becomes synergy minus the redundancy *counted multiple times*.

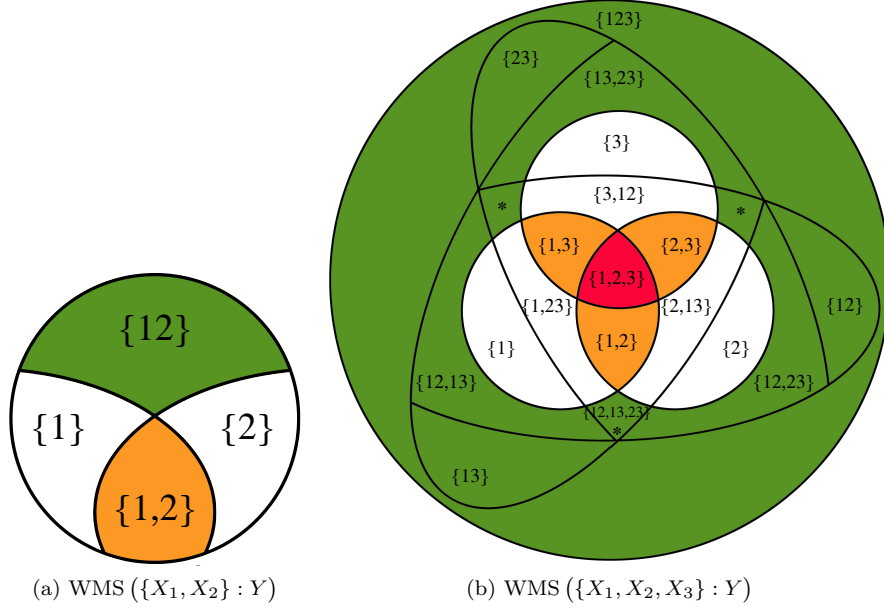


Figure 2.2: PI-diagrams illustrating WholeMinusSum synergy for  $n = 2$  (left) and  $n = 3$  (right). The colors denote the added and subtracted PI-regions.  $\text{WMS}(\mathbf{X} : Y)$  is the green PI-region(s) minus the orange PI-region(s) minus two times any red PI-region.

#### 2.1.4 WholeMinusPartitionSum: $\text{WMPS}(\mathbf{X} : Y)$

WholeMinusPartitionSum, denoted  $\text{WMPS}(\mathbf{X} : Y)$ , is a stricter generalization of WMS synergy for  $n > 2$ . It was introduced in [34, 1] and is defined as,

$$\text{WMPS}(\mathbf{X} : Y) \equiv I(X_{1\dots n} : Y) - \max_{\mathbf{P}} \sum_{i=1}^{|\mathbf{P}|} I(P_i : Y) , \quad (2.6)$$

where  $\mathbf{P}$  enumerates over all partitions of the set of predictors  $\{X_1, \dots, X_n\}$ . WholeMinusPartitionSum is a signed measure where a positive value signifies synergy and a negative value signifies redundancy.

For  $n = 3$ , there are four partitions of  $\mathbf{X}$  resulting in four possible PI-diagrams—one for each partition. Figure 2.3 depicts the four possible values of  $\text{WMPS}(\{X_1, X_2, X_3\} : Y)$ . Because  $\{X_1, \dots, X_n\}$  is a possible partition of  $\mathbf{X}$ ,  $\text{WMPS}(\mathbf{X} : Y) \leq \text{WMS}(\mathbf{X} : Y)$ .

<sup>3</sup>This is deeper than [29]’s point that a mish-mash of synergy and redundancy across different states of  $y \in Y$  can average to zero. E.g., Figure 2.6 evaluates to zero for *every state*  $y \in Y$ .

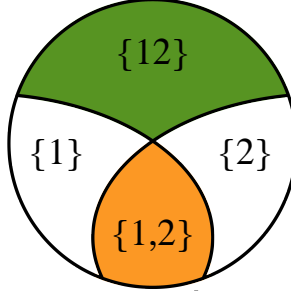
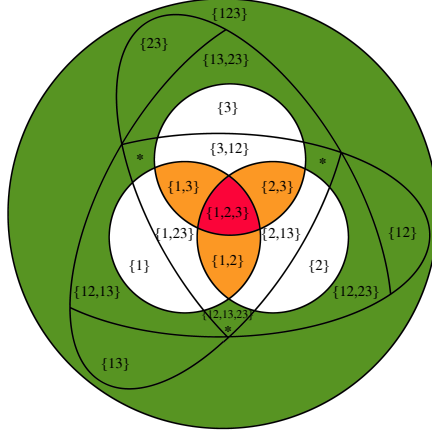
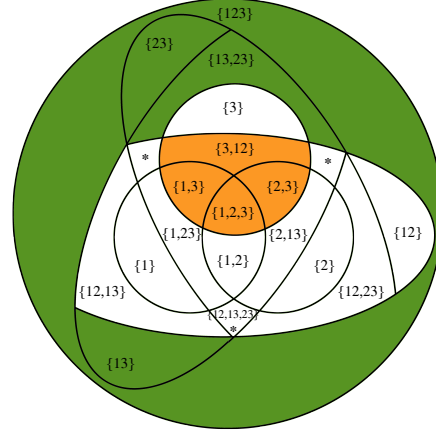
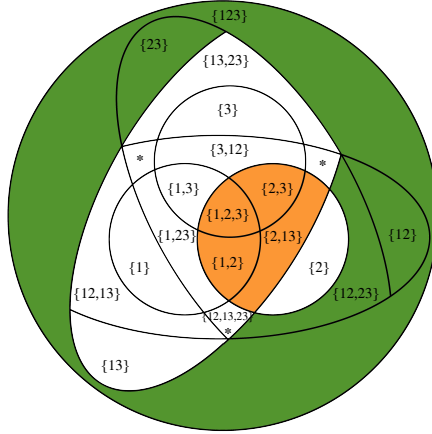
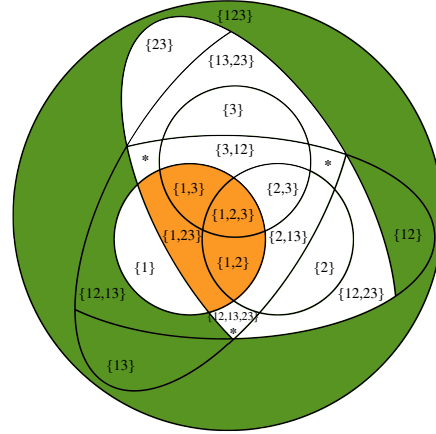
(a)  $\text{W MPS}(\{X_1, X_2\} : Y)$ (b)  $\mathbf{P} = \{X_1, X_2, X_3\}$ (c)  $\mathbf{P} = \{X_1 X_2, X_3\}$ (d)  $\mathbf{P} = \{X_1 X_3, X_2\}$ (e)  $\mathbf{P} = \{X_2 X_3, X_1\}$ 

Figure 2.3: PI-diagrams depicting WholeMinusPartitionSum synergy for  $n = 2$  (2.3a) and  $n = 3$  (2.3b–2.3e). Each measure is the green PI-regions minus the orange PI-regions minus two times any red PI-region.  $\text{W MPS}(\{X_1, X_2, X_3\} : Y)$  is the *minimum* value over subfigures 2.3b–2.3e.

### 2.1.5 $I_{\max}$ synergy: $\mathcal{S}_{\max}(\mathbf{X} : Y)$

$I_{\max}$  synergy, denoted  $\mathcal{S}_{\max}$ , was the first synergy measure derived from Partial Information Decomposition[36].

$\mathcal{S}_{\max}$  defines synergy as the whole beyond the state-dependent *maximum* of its elements,

$$\mathcal{S}_{\max}(\mathbf{X} : Y) \equiv I(X_{1\dots n} : Y) - I_{\max}(\{X_1, \dots, X_n\} : Y) \quad (2.7)$$

$$= I(X_{1\dots n} : Y) - \sum_{y \in Y} \Pr(Y = y) \max_i I(X_i : Y = y) , \quad (2.8)$$

where  $I(X_i : Y = y)$  is [10]’s “specific-surprise”,

$$I(X_i : Y = y) \equiv D_{\text{KL}} \left[ \Pr(X_i | y) \parallel \Pr(X_i) \right] \quad (2.9)$$

$$= \sum_{x_i \in X_i} \Pr(x_i | y) \log \frac{\Pr(x_i, y)}{\Pr(x_i) \Pr(y)} . \quad (2.10)$$

There are two major advantages of  $\mathcal{S}_{\max}$  synergy.  $\mathcal{S}_{\max}$  is not only nonnegative, but also invariant to duplicate predictors.

### 2.1.6 Correlational importance: $\Delta I(\mathbf{X}; Y)$

Correlational importance, denoted  $\Delta I$ , comes from [27, 25, 26, 28, 21]. Correlational importance quantifies the “informational importance of conditional dependence” or the “information lost when ignoring conditional dependence” among the predictors decoding target  $Y$ . On casual inspection,  $\Delta I$  seems related to our intuitive conception of synergy.  $\Delta I$  is defined as,

$$\Delta I(\mathbf{X}; Y) \equiv D_{\text{KL}} \left[ \Pr(Y | X_{1\dots n}) \parallel \Pr_{\text{ind}}(Y | \mathbf{X}) \right] \quad (2.11)$$

$$= \sum_{y, \mathbf{x} \in Y, \mathbf{X}} \Pr(y, x_{1\dots n}) \log \frac{\Pr(y | x_{1\dots n})}{\Pr_{\text{ind}}(y | \mathbf{x})} , \quad (2.12)$$

where  $\Pr_{\text{ind}}(y | \mathbf{x}) \equiv \frac{\Pr(y) \prod_{i=1}^n \Pr(x_i | y)}{\sum_{y'} \Pr(y') \prod_{i=1}^n \Pr(x_i | y')}$ . After some algebra<sup>4</sup> eq. (2.12) becomes,

$$\Delta I(\mathbf{X}; Y) = \text{TC}(X_1; \dots; X_n | Y) - D_{\text{KL}} \left[ \Pr(X_{1\dots n}) \parallel \sum_y \Pr(y) \prod_{i=1}^n \Pr(X_i | y) \right] . \quad (2.13)$$

$\Delta I$  is conceptually innovative, yet examples reveal that  $\Delta I$  measures something ever-so-subtly different from intuitive synergistic information.

---

<sup>4</sup>See Appendix 2.A for the steps between eqs. (2.12) and (2.13).

## 2.2 The six prior measures are not equivalent

For  $n = 2$ , the four measures  $\mathcal{S}_{\text{MMI}}$ ,  $\mathcal{II}$ , WMS, and WMPS are equivalent. But in general, none of these six measures are equivalent. Example AND (Figure 2.4) shows that  $\mathcal{S}_{\text{max}}$  and  $\Delta I$  are not equivalent. Example XORMULTICOAL (Figure 2.5) shows that  $\mathcal{S}_{\text{MMI}}$ ,  $\mathcal{II}$ , WMS, and WMPS are not equivalent.

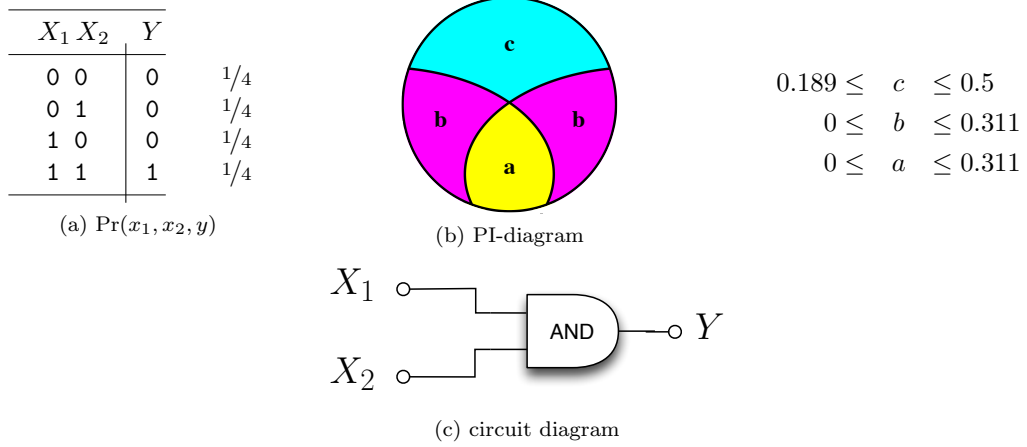


Figure 2.4: Example AND. The exact PI-decomposition of an AND-gate remains uncertain. But we can bound  $a$ ,  $b$ , and  $c$  using WMS and  $\mathcal{S}_{\text{max}}$ .

Example	$\mathcal{S}_{\text{MMI}}$	$\mathcal{II}$	WMS	WMPS	$\mathcal{S}_{\text{max}}$	$\Delta I$
AND	0.189	0.189	0.189	0.189	$1/2$	0.104
XORMULTICOAL	2	-2	1	0	1	1

Table 2.1: Examples demonstrating that the six prior measures are not equivalent.

## 2.3 Counter-intuitive behaviors of the six prior measures

### 2.3.1 $\mathcal{I}_{\text{max}}$ synergy: $\mathcal{S}_{\text{max}}$

Despite several desired properties,  $\mathcal{S}_{\text{max}}$  sometimes miscategorizes merely unique information as synergistic. This can be seen in example UNQ (Figure 1.4). In example UNQ, the wires in Figure 1.4b don't even touch, yet  $\mathcal{S}_{\text{max}}$  asserts there is one bit of synergy and one bit of redundancy—this is palpably strange.

A more abstract way to understand why  $\mathcal{S}_{\text{max}}$  overestimates synergy is to imagine a hypothetical example where there are exactly two bits of unique information for every state  $y \in Y$  and no synergy or redundancy.  $\mathcal{S}_{\text{max}}$  would be the whole (both unique bits) minus the *maximum* over both



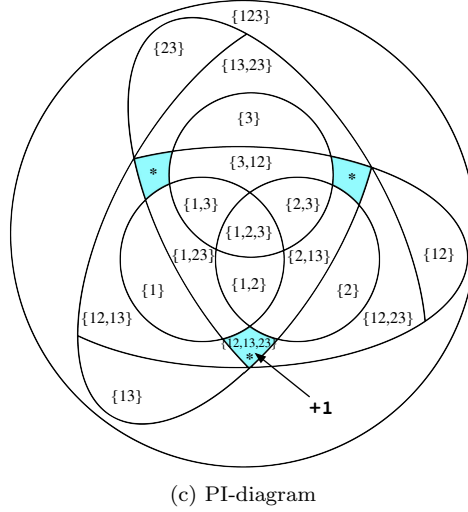
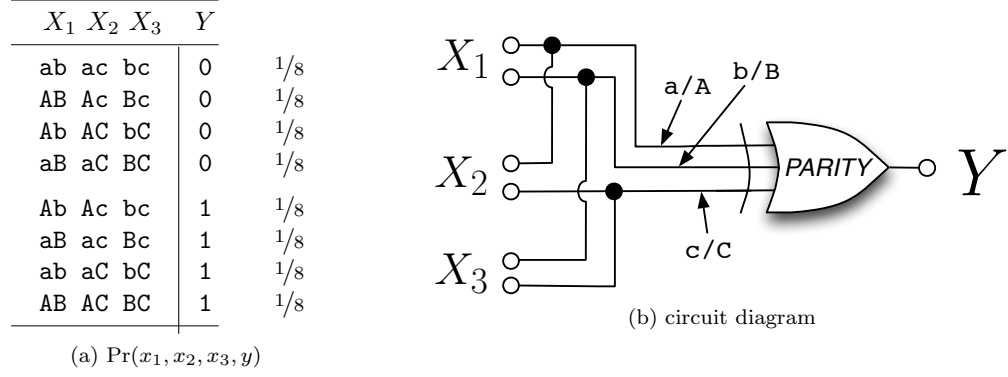


Figure 2.5: Example XORMULTICOAL demonstrates how the same information can be specified by multiple coalitions. In XORMULTICOAL the target  $Y$  has one bit of uncertainty,  $H(Y) = 1$  bit, and  $Y$  is the *parity* of three incoming wires. Just as the output of XOR is specified only after knowing the state of both inputs, the output of XORMULTICOAL is specified only after knowing the state of all three wires. Each predictor is distinct and has access to two of the three incoming wires. For example, predictor  $X_1$  has access to the  $a/A$  and  $b/B$  wires,  $X_2$  has access to the  $a/A$  and  $c/C$  wires, and  $X_3$  has access to the  $b/B$  and  $c/C$  wires. Although no single predictor specifies  $Y$ , any coalition of two predictors has access to all three wires and fully specifies  $Y$ ,  $I(X_1X_2:Y) = I(X_1X_3:Y) = I(X_2X_3:Y) = H(Y) = 1$  bit. In the PI-diagram this puts one bit in PI-region  $\{12, 13, 23\}$  and zero everywhere else.

predictors, which would be the  $\max[1, 1] = 1$  bit. The  $S_{\max}$  synergy would then be  $2 - 1 = 1$  bit of synergy, even though by definition there was no synergy, but merely two bits of unique information.

Altogether, we conclude that  $S_{\max}$  *overestimates* the intuitive synergy by miscategorizing merely unique information as synergistic whenever two or more predictors have unique information about the target.

### 2.3.2 $S_{\text{MMI}}$ , $II$ , WMS, WMPS

All four of these measures are equivalent for  $n = 2$ . Given this agreement, it is ironic that there are counter-intuitive examples even for  $n = 2$ . A concrete example demonstrating a “synergy minus

redundancy” behavior for  $n = 2$  is example RDNXOR (Figure 2.6), which overlays examples RDN and XOR to form a single system. The target  $Y$  has two bits of uncertainty, i.e.  $H(Y) = 2$ . Like RDN, either  $X_1$  or  $X_2$  identically specifies the letter of  $Y$  ( $\mathbf{r}/\mathbf{R}$ ), making one bit of redundant information. Like XOR, only the coalition  $X_1X_2$  specifies the digit of  $Y$  ( $\mathbf{0}/\mathbf{1}$ ), making one bit of synergistic information. Together this makes one bit of redundancy and one bit of synergy. We assert that for  $n = 2$ , all four measures *underestimate* the synergy. Equivalently, we say that their answer for  $n = 2$  is a *lowerbound* on the intuitive synergy.

Note that in RDNXOR every state  $y \in Y$  conveys one bit of redundant information and one bit of synergistic information, e.g. for the state  $y = \mathbf{r0}$  the letter “ $\mathbf{r}$ ” is specified redundantly and the digit “ $\mathbf{0}$ ” is specified synergistically.

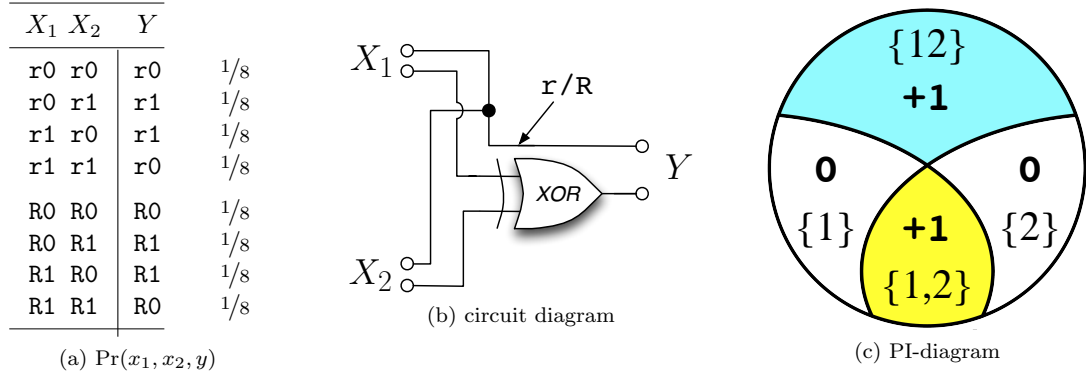


Figure 2.6: Example RDNXOR has one bit of redundancy and one bit of synergy. Yet for this example, the four most common measures of synergy arrive at zero bits.

Our next example, PARITYRDN (Figure 2.7), has one bit of synergy and two bits of redundancy for a total of  $I(X_1X_2X_3:Y) = H(Y) = 3$  bits. It emphasizes the disagreement between  $\mathcal{II}$  and measures  $\mathcal{S}_{\text{MMI}}$ , WMS, and WMPS. If  $\mathcal{S}_{\text{MMI}}$ , WMS, or WMPS were always simply “synergy minus redundancy”, then one of them would calculate  $1 - 2 = -1$  bits. But for this example all three measures subtracts redundancies *multiple times* to calculate  $1 - (2 \cdot 2) = -3$  bits, signifying all three bits of  $H(Y)$  are specified redundantly.  $\mathcal{II}$  makes a different misstep. Instead of subtracting redundancy multiple times, for  $n = 3$   $\mathcal{II}$  *adds* the maximum redundancy to calculate  $1 + 2 = +3$  bits, signifying three bits of synergy and no redundancy. Both answers are palpably mistaken.

### 2.3.3 Correlational importance: $\Delta I$

The first concerning example is [29]’s Figure 4, where  $\Delta I$  exceeds the mutual information  $I(X_{1..n}:Y)$  with  $\Delta I(\mathbf{X};Y) = 0.0145$  and  $I(X_{1..n}:Y) = 0.0140$ . This fact alone prevents interpreting  $\Delta I$  the magnitude of mutual  $I(X_{1..n}:Y)$  arising from correlational dependence.

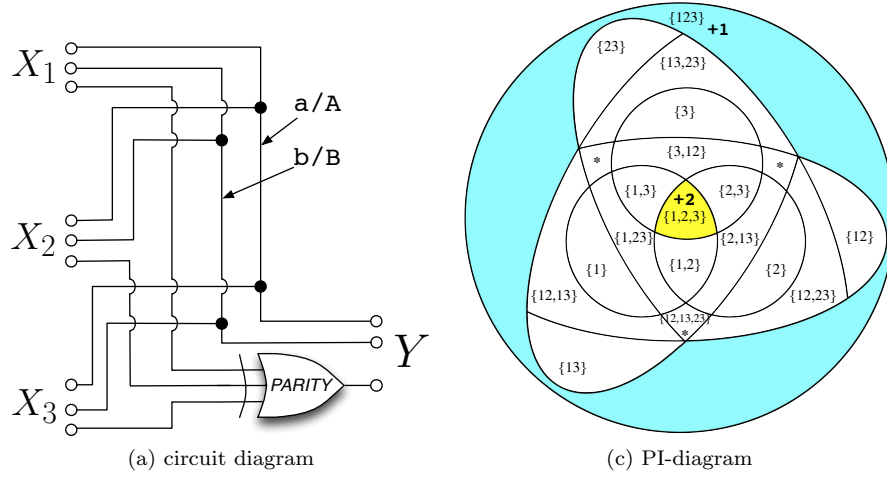
Could  $\Delta I$  upperbound synergy instead? We turn to example AND (Figure 2.4) with  $n = 2$  independent binary predictors and target  $Y$  is the AND of  $X_1$  and  $X_2$ . Although AND’s exact

PI-region decomposition remains uncertain, we can still bound the synergy. For example, the  $\text{WMS}(\{X_1, X_2\} : Y) \approx 0.189$  and  $\mathcal{S}_{\max}(\{X_1, X_2\} : Y) = 0.5$  bits. So we know the synergy must be between  $(0.189, 0.5]$  bits. Despite this,  $\Delta I(\mathbf{X}; Y) = 0.104$  bits, thus  $\Delta I$  does not upperbound synergy either.

Taking both together, we conclude that  $\Delta I$  measures something fundamentally different from synergistic information.

Example	$\mathcal{S}_{\text{MMI}}$	$\mathcal{II}$	WMS	WMPS	$\mathcal{S}_{\max}$	$\Delta I$
UNQ	0	0	0	0	1	0
RDNXOR	0	0	0	0	1	1
PARITYRDN	-3	3	-3	-3	1	1
AND	0.189	0.189	0.189	0.189	1/2	0.104

Table 2.2: Examples demonstrating that all six prior measures have shortcomings.



$X_1$	$X_2$	$X_3$	$Y$		$X_1$	$X_2$	$X_3$	$Y$	
ab0	ab0	ab0	ab0	$1/32$	Ab0	Ab0	Ab0	Ab0	$1/32$
ab0	ab1	ab1	ab0	$1/32$	Ab0	Ab1	Ab1	Ab0	$1/32$
ab1	ab0	ab1	ab0	$1/32$	Ab1	Ab0	Ab1	Ab0	$1/32$
ab1	ab1	ab0	ab0	$1/32$	Ab1	Ab1	Ab0	Ab0	$1/32$
aB0	aB0	aB0	aB0	$1/32$	AB0	AB0	AB0	AB0	$1/32$
aB0	aB1	aB1	aB0	$1/32$	AB0	AB1	AB1	AB0	$1/32$
aB1	aB0	aB1	aB0	$1/32$	AB1	AB0	AB1	AB0	$1/32$
aB1	aB1	aB0	aB0	$1/32$	AB1	AB1	AB0	AB0	$1/32$
ab0	ab0	ab1	ab1	$1/32$	Ab0	Ab0	Ab1	Ab1	$1/32$
ab0	ab1	ab0	ab1	$1/32$	Ab0	Ab1	Ab0	Ab1	$1/32$
ab1	ab0	ab0	ab1	$1/32$	Ab1	Ab0	Ab0	Ab1	$1/32$
ab1	ab1	ab1	ab1	$1/32$	Ab1	Ab1	Ab1	Ab1	$1/32$
aB0	aB0	aB1	aB1	$1/32$	AB0	AB0	AB1	AB1	$1/32$
aB0	aB1	aB0	aB1	$1/32$	AB0	AB1	AB0	AB1	$1/32$
aB1	aB0	aB0	aB1	$1/32$	AB1	AB0	AB0	AB1	$1/32$
aB1	aB1	aB1	aB1	$1/32$	AB1	AB1	AB1	AB1	$1/32$

(b)  $\Pr(x_1, x_2, x_3, y)$ 

Figure 2.7: Example PARITYRDN. Three predictors redundantly specify two bits of  $Y$ ,  $I(X_1:Y) = I(X_2:Y) = I(X_3:Y) = 2$  bits. At the same time, the three predictors holistically specify the third and final bit of  $Y$ ,  $I(X_1X_2X_3:Y) = H(Y) = 3$  bits.

# Appendix

## 2.A Algebraic simplification of $\Delta I$

Prior literature [25, 26, 28, 21] defines  $\Delta I(\mathbf{X}; Y)$  as,

$$\Delta I(\mathbf{X}; Y) \equiv D_{\text{KL}} \left[ \Pr(Y|X_{1\dots n}) \parallel \Pr_{\text{ind}}(Y|\mathbf{X}) \right] \quad (2.14)$$

$$= \sum_{\mathbf{x}, y \in \mathbf{X}, Y} \Pr(\mathbf{x}, y) \log \frac{\Pr(y|\mathbf{x})}{\Pr_{\text{ind}}(y|\mathbf{x})}. \quad (2.15)$$

Where,

$$\Pr_{\text{ind}}(Y = y|\mathbf{X} = \mathbf{x}) \equiv \frac{\Pr(y) \Pr_{\text{ind}}(\mathbf{X} = \mathbf{x}|Y = y)}{\Pr_{\text{ind}}(\mathbf{X} = \mathbf{x})} \quad (2.16)$$

$$= \frac{\Pr(y) \prod_{i=1}^n \Pr(x_i|y)}{\Pr_{\text{ind}}(\mathbf{x})} \quad (2.17)$$

$$\Pr_{\text{ind}}(\mathbf{X} = \mathbf{x}) \equiv \sum_{y \in Y} \Pr(Y = y) \prod_{i=1}^n \Pr(x_i|y) \quad (2.18)$$

The definition of  $\Delta I$ , eq. (2.14), reduces to,

$$\Delta I(\mathbf{X}; Y) = \sum_{\mathbf{x}, y \in \mathbf{X}, Y} \Pr(\mathbf{x}, y) \log \frac{\Pr(y|\mathbf{x})}{\Pr_{\text{ind}}(y|\mathbf{x})} \quad (2.19)$$

$$= \sum_{\mathbf{x}, y \in \mathbf{X}, Y} \Pr(\mathbf{x}, y) \log \frac{\Pr(y|\mathbf{x}) \Pr_{\text{ind}}(\mathbf{x})}{\Pr(y) \prod_{i=1}^n \Pr(x_i|y)} \quad (2.20)$$

$$= \sum_{\mathbf{x}, y \in \mathbf{X}, Y} \Pr(\mathbf{x}, y) \log \frac{\Pr(\mathbf{x}|y)}{\prod_{i=1}^n \Pr(x_i|y)} \frac{\Pr_{\text{ind}}(\mathbf{x})}{\Pr(\mathbf{x})} \quad (2.21)$$

$$= \sum_{\mathbf{x}, y \in \mathbf{X}, Y} \Pr(\mathbf{x}, y) \log \frac{\Pr(\mathbf{x}|y)}{\prod_{i=1}^n \Pr(x_i|y)} + \sum_{\mathbf{x}, y \in \mathbf{X}, Y} \Pr(\mathbf{x}, y) \log \frac{\Pr_{\text{ind}}(\mathbf{x})}{\Pr(\mathbf{x})} \quad (2.22)$$

$$= \sum_{\mathbf{x}, y \in \mathbf{X}, Y} \Pr(\mathbf{x}, y) \log \frac{\Pr(\mathbf{x}|y)}{\prod_{i=1}^n \Pr(x_i|y)} - \sum_{\mathbf{x} \in \mathbf{X}} \Pr(\mathbf{x}) \log \frac{\Pr(\mathbf{x})}{\Pr_{\text{ind}}(\mathbf{x})} \quad (2.23)$$

$$= \text{D}_{\text{KL}} \left[ \Pr(X_{1\dots n}|Y) \left\| \prod_{i=1}^n \Pr(X_i|Y) \right\| \right] - \text{D}_{\text{KL}} [\Pr(X_{1\dots n}) \| \Pr_{\text{ind}}(\mathbf{X})]$$

$$= \text{TC}(X_1; \dots; X_n|Y) - \text{D}_{\text{KL}} [\Pr(X_{1\dots n}) \| \Pr_{\text{ind}}(\mathbf{X})] .$$

where  $\text{TC}(X_1; \dots; X_n|Y)$  is the conditional total correlation among the predictors given  $Y$ .

## Part II

# Making Progress

## Chapter 3

# First Nontrivial Lowerbound on Synergy

**Remark:** This chapter borrows liberally from the joint paper [14].

### 3.1 Introduction

Introduced in [36], *Partial Information Decomposition* (PID) is an immensely useful framework for deepening our understanding of multivariate interactions, particularly our understanding of informational redundancy and synergy. To harness the PID framework, the user brings her own measure of *intersection information*,  $I_{\cap}(X_1, \dots, X_n : Y)$ , which quantifies the magnitude of information that each of the  $n$  predictors  $X_1, \dots, X_n$  conveys “redundantly” about a target random variable  $Y$ . An antichain lattice of redundant, unique, and synergistic partial informations is built from the intersection information.

In [36], the authors propose to use the following quantity,  $I_{\min}$ , as the intersection information measure:

$$\begin{aligned} I_{\min}(X_1, \dots, X_n : Y) &\equiv \sum_y \Pr(y) \min_i I(X_i : Y = y) \\ &= \sum_y \Pr(y) \min_i D_{\text{KL}} \left[ \Pr(X_i | y) \parallel \Pr(X_i) \right], \end{aligned} \quad (3.1)$$

where  $D_{\text{KL}}$  is the Kullback-Leibler divergence.

Though  $I_{\min}$  is an intuitive and plausible choice for the intersection information, [15] showed that  $I_{\min}$  has counterintuitive properties. In particular,  $I_{\min}$  calculates one bit of redundant information for example UNQ (Figure 3.3). It does this because each input shares one bit of information with the output. However, it is quite clear that the shared informations are, in fact, different:  $X_1$  provides the low bit, while  $X_2$  provides the high bit. This led to the conclusion that  $I_{\min}$  *over-estimates* the ideal intersection information measure by focusing only on *how much* information the inputs



provide to the output. An ideal measure of intersection information must recognize that there are non-equivalent ways of providing information to the output. The search for an improved intersection information measure ensued, continued through [18, 7, 23], and today a widely accepted intersection information measure remains undiscovered.

Here we do not definitively solve this problem, but we present a strong candidate intersection information measure for the special case of *zero-error* information. This is useful in of itself because it provides a template for how the yet undiscovered ideal intersection information measure for Shannon mutual information could work. Alternatively, if a Shannon intersection information measure with the same properties does not exist, then we have learned something significant.

In the next section, we introduce some definitions, some notation, and a necessary lemma. We also extend and clarify the desired properties for intersection information. In Section 3.4 we introduce zero-error information and its intersection information measure. In Section 3.5 we use the same methodology to produce a novel candidate for the Shannon intersection information. In Section 3.6 we show the successes and shortcomings of our candidate intersection information measure using example circuits. Finally, in Section 3.8 we summarize our progress towards the ideal intersection information measure and suggest directions for improvement. The Appendix is devoted to technical lemmas and their proofs, to which we refer in the main text.

## 3.2 Two examples elucidating desired properties for synergy

To help the reader develop intuition for any proper measure of synergy, we illustrate some desired properties of synergistic information with pedagogical examples. Both examples are derived from example XOR.

### 3.2.1 XorDuplicate: Synergy is invariant to duplicating a predictor

Example XORDUPLICATE (Figure 3.1) adds a third predictor,  $X_3$ , a copy of predictor  $X_1$ , to XOR. Whereas in XOR the target  $Y$  is specified only by coalition  $X_1X_2$ , duplicating predictor  $X_1$  as  $X_3$  makes the target equally specifiable by coalition  $X_3X_2$ .

Although now two different coalitions identically specify  $Y$ , mutual information is invariant to duplicates, e.g.  $I(X_1X_2X_3:Y) = I(X_1X_2:Y)$  bit. For synergistic information to be likewise bounded between zero and the total mutual information  $I(X_{1..n}:Y)$ , synergistic information must similarly be invariant to duplicates, e.g. the synergistic information between set  $\{X_1, X_2\}$  and  $Y$  must be the same as the synergistic information between  $\{X_1, X_2, X_3\}$  and  $Y$ . This makes sense because if synergistic information is defined as the information in the whole beyond its parts, duplicating a part does not increase the net information provided by the parts. Altogether, we assert that *duplicating a predictor does not change the synergistic information*.

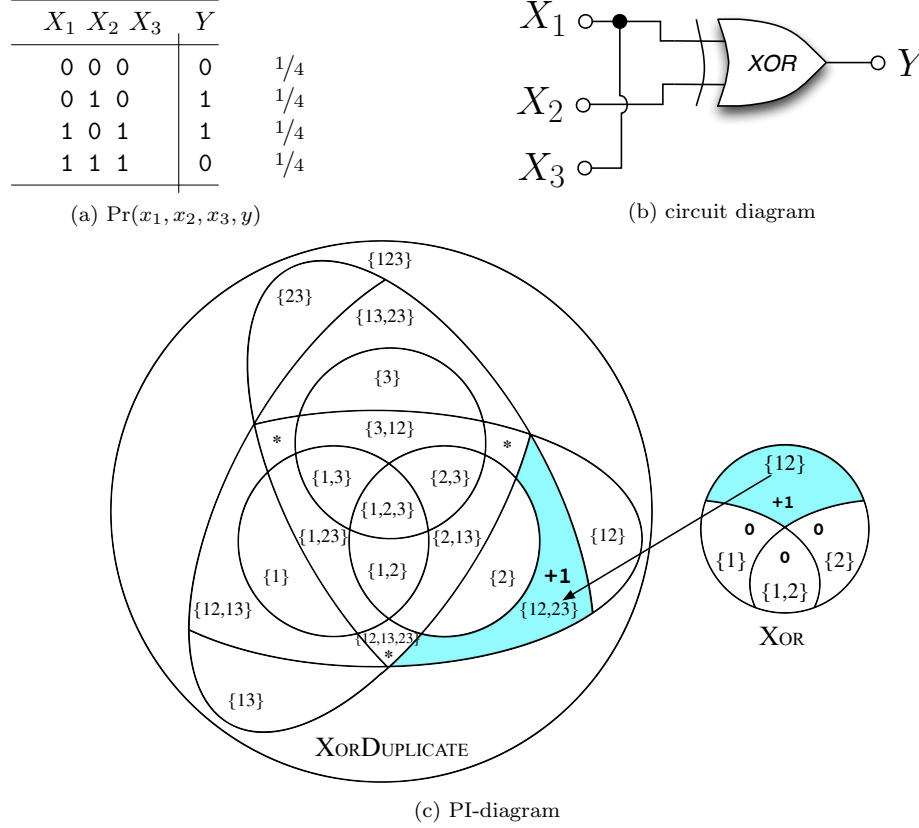


Figure 3.1: Example XORDUPLICATE shows that duplicating predictor  $X_1$  as  $X_3$  turns the single-coalition synergy  $\{12\}$  into the multi-coalition synergy  $\{12, 23\}$ . After duplicating  $X_1$ , the coalition  $X_3X_2$  as well as coalition  $X_1X_2$  specifies  $Y$ . Synergistic information is unchanged from XOR,  $I(X_3X_2:Y) = I(X_1X_2:Y) = H(Y) = 1$  bit.

### 3.2.2 XorLoses: Adding a new predictor can decrease synergy

Example XORLOSES (Figure 3.2) adds a third predictor,  $X_3$ , to XOR and concretizes the distinction between synergy and “redundant synergy”. In XORLOSES the target  $Y$  has one bit of uncertainty, and just as in example XOR the coalition  $X_1X_2$  fully specifies the target,  $I(X_1X_2:Y) = H(Y) = 1$  bit. However, XORLOSES has *zero* intuitive synergy because the newly added singleton predictor,  $X_3$ , fully specifies  $Y$  by itself. This makes the synergy between  $X_1$  and  $X_2$  *completely redundant*—everything the coalition  $X_1X_2$  specifies is now already specified by the singleton  $X_3$ .

## 3.3 Preliminaries

### 3.3.1 Informational Partial Order and Equivalence

We assume an underlying probability space on which we define random variables, as denoted by capital letters (e.g.,  $X$ ,  $Y$ , and  $Z$ ). In this paper, we consider only random variables taking values

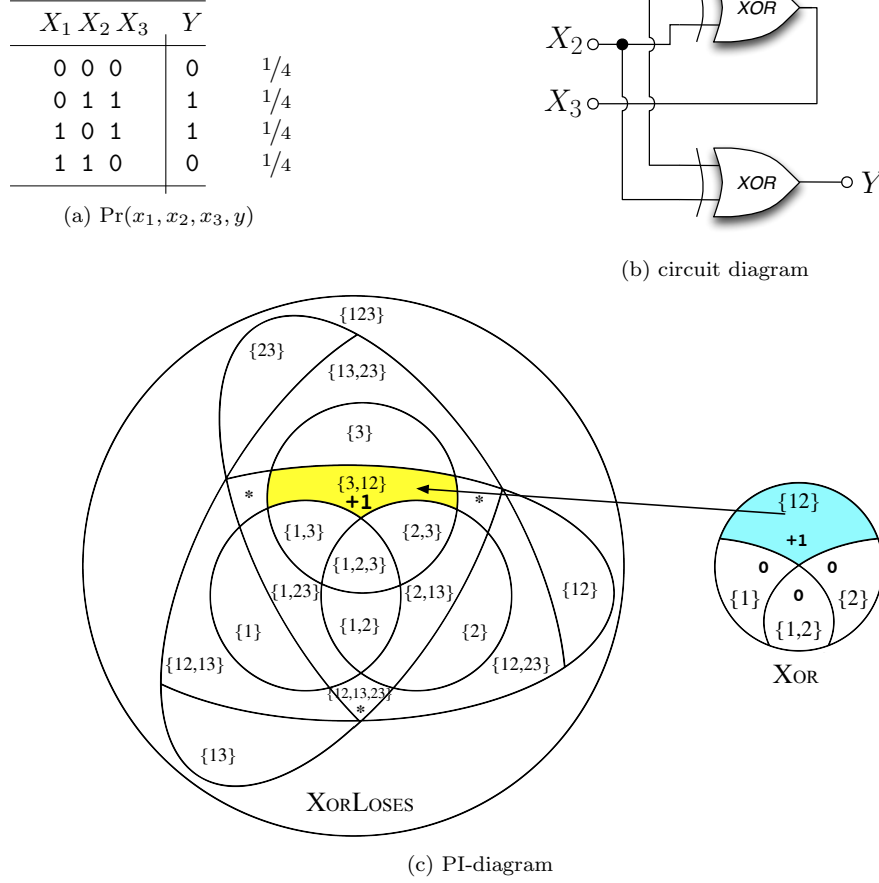


Figure 3.2: Example XORLOSES. Target  $Y$  is fully specified by the coalition  $X_1X_2$  as well as by the singleton  $X_3$ .  $I(X_1X_2:Y) = I(X_3:Y) = H(Y) = 1$  bit. Therefore, the information synergistically specified by coalition  $X_1X_2$  is a redundant synergy.

on finite spaces.

Given random variables  $X$  and  $Y$ , we write  $X \preceq Y$  to signify that there exists a measurable function  $f$  such that  $X = f(Y)$ . In this case, following the terminology in [22], we say that  $X$  is *informationally poorer* than  $Y$ ; this induces a partial order on the set of random variables. Similarly, we write  $X \succeq Y$  if  $Y \preceq X$ , in which case we say  $X$  is *informationally richer* than  $Y$ .

If  $X$  and  $Y$  are such that  $X \preceq Y$  and  $X \succeq Y$ , then we write  $X \cong Y$ . In this case, again following [22], we say that  $X$  and  $Y$  are *informationally equivalent*. In other words,  $X \cong Y$  if and only if there's an invertible function between  $X$  and  $Y$ , i.e., one can relabel the values of  $X$  to obtain a random value that is equal to  $Y$  and vice versa.

This “information-equivalence” relation can easily be shown to be an equivalence relation, so that we can partition the set of all random variables into disjoint equivalence classes. The  $\preceq$  ordering is invariant within these equivalence classes in the following sense: if  $X \preceq Y$  and  $Y \cong Z$ , then  $X \preceq Z$ . Similarly, if  $X \preceq Y$  and  $X \cong Z$ , then  $Z \preceq Y$ . Moreover, within each equivalence class, the entropy

is invariant, as stated formally in Lemma 3.3.1 below.

### 3.3.2 Information Lattice

Next, we follow [22] and consider the *join* and *meet* operators. These operators were defined for *information elements*, which are  $\sigma$ -algebras, or, equivalently, equivalence classes of random variables. We deviate from [22], though, by defining the join and meet operators for random variables, but we do preserve their conceptual properties.

Given random variables  $X$  and  $Y$ , we define  $X \vee Y$  (called the *join* of  $X$  and  $Y$ ) to be an informationally poorest (“smallest” in the sense of the partial order  $\preceq$ ) random variable such that  $X \preceq X \vee Y$  and  $Y \preceq X \vee Y$ . In other words, if  $Z$  is such that  $X \preceq Z$  and  $Y \preceq Z$ , then  $X \vee Y \preceq Z$ . Note that  $X \vee Y$  is unique only up to equivalence with respect to  $\cong$ . In other words,  $X \vee Y$  does not define a specific, unique random variable. Nonetheless, standard information-theoretic quantities are invariant over the set of random variables satisfying the condition specified above. For example, the entropy of  $X \vee Y$  is invariant over the entire equivalence class of random variables satisfying the condition above (by Lemma 3.3.1(a) below). Similarly, the inequality  $Z \preceq X \vee Y$  does not depend on the specific random variable chosen, as long as it satisfies the condition above. Note that the pair  $(X, Y)$  is an instance of  $X \vee Y$ .

In a similar vein, given random variables  $X$  and  $Y$ , we define  $X \wedge Y$  (called the *meet* of  $X$  and  $Y$ ) to be an informationally richest random variable (“largest” in the sense of  $\succeq$ ) such that  $X \wedge Y \preceq X$  and  $X \wedge Y \preceq Y$ . In other words, if  $Z$  is such that  $Z \preceq X$  and  $Z \preceq Y$ , then  $Z \preceq X \wedge Y$ . Following [16], we also call  $X \wedge Y$  the *common random variable* of  $X$  and  $Y$ . Again, considering the entropy of  $X \wedge Y$  or the inequality  $Z \preceq X \wedge Y$  does not depend on the specific random variable chosen, as long as it satisfies the condition above.

The  $\vee$  and  $\wedge$  operators satisfy the algebraic properties of a *lattice* [22]. In particular, the following hold:

- commutative laws:  $X \vee Y \cong Y \vee X$  and  $X \wedge Y \cong Y \wedge X$
- associative laws:  $X \vee (Y \vee Z) \cong (X \vee Y) \vee Z$  and  $X \wedge (Y \wedge Z) \cong (X \wedge Y) \wedge Z$
- absorption laws:  $X \vee (X \wedge Y) \cong X$  and  $X \wedge (X \vee Y) \cong X$
- idempotent laws:  $X \vee X \cong X$  and  $X \wedge X \cong X$
- generalized absorption laws: if  $X \preceq Y$ , then  $X \vee Y \cong Y$  and  $X \wedge Y \cong X$  .

Finally, the partial order  $\preceq$  is preserved under  $\vee$  and  $\wedge$ , i.e., if  $X \preceq Y$ , then  $X \vee Z \preceq Y \vee Z$  and  $X \wedge Z \preceq Y \wedge Z$ .

### 3.3.3 Invariance and Monotonicity of Entropy

Let  $H(\cdot)$  represent the entropy function, and  $H(\cdot|\cdot)$  the conditional entropy. To be consistent with the colon in the intersection information, we denote the Shannon mutual information between  $X$  and  $Y$  by  $I(X:Y)$  instead of the more common  $I(X;Y)$ . Lemma 3.3.1 establishes the invariance and monotonicity of the entropy and conditional entropy functions with respect to  $\cong$  and  $\preceq$ .

**Lemma 3.3.1.** *The following hold:*

- (a) *If  $X \cong Y$ , then  $H(X) = H(Y)$ ,  $H(X|Z) = H(Y|Z)$ , and  $H(Z|X) = H(Z|Y)$ .*
- (b) *If  $X \preceq Y$ , then  $H(X) \leq H(Y)$ ,  $H(X|Z) \leq H(Y|Z)$ , and  $H(Z|X) \geq H(Z|Y)$ .*
- (c)  *$X \preceq Y$  if and only if  $H(X|Y) = 0$ .*

*Proof.* Part (a) follows from [22], Proposition 1. Part (c) follows from [22], Proposition 4. The first two desired inequalities in part (b) follow from [22], Proposition 5. Now we show that if  $X \preceq Y$ , then  $H(Z|X) \geq H(Z|Y)$ . Suppose that  $X \preceq Y$ . Then, by the generalized absorption law,  $X \vee Y \cong Y$ . We have

$$\begin{aligned}
 I(Z:Y) &= H(Y) - H(Y|Z) \\
 &= H(X \vee Y) - H(X \vee Y|Z) \quad \text{by part (a)} \\
 &= I(Z:X \vee Y) \\
 &= I(Z:X) + I(Z:Y|X) \\
 &\geq I(Z:X).
 \end{aligned}$$

Substituting  $I(Z:Y) = H(Z) - H(Z|Y)$  and  $I(Z:X) = H(Z) - H(Z|X)$ , we obtain  $H(Z|X) \geq H(Z|Y)$  as desired.  $\square$

**Remark:** Because  $(X, Y) \cong X \vee Y$  as noted before, we also have  $H(X, Y) = H(X \vee Y)$  by Lemma 3.3.1(a).

### 3.3.4 Desired Properties of Intersection Information

There are currently 12 intuitive properties that we wish the ideal intersection information measure  $I_\cap$  to satisfy. Some are new (e.g. **(M<sub>1</sub>)**, **(Eq)**, **(LB)**), but most were introduced earlier, in various forms, Refs. [36, 15, 18, 7, 23]. They are as follows:

**(GP)** Global Positivity:  $I_\cap(X_1, \dots, X_n:Y) \geq 0$ , and  $I_\cap(X_1, \dots, X_n:Y) = 0$  if  $Y$  is a constant.

**(Eq)** Equivalence-Class Invariance:  $I_\cap(X_1, \dots, X_n:Y)$  is invariant under substitution of  $X_i$  (for any  $i = 1, \dots, n$ ) or  $Y$  by an informationally equivalent random variable.

(**TM**) Target Monotonicity: If  $Y \preceq Z$ , then  $I_{\cap}(X_1, \dots, X_n : Y) \leq I_{\cap}(X_1, \dots, X_n : Z)$ .

(**M<sub>0</sub>**) Weak Monotonicity:  $I_{\cap}(X_1, \dots, X_n, W : Y) \leq I_{\cap}(X_1, \dots, X_n : Y)$  with equality if there exists  $Z \in \{X_1, \dots, X_n\}$  such that  $Z \preceq W$ .

(**S<sub>0</sub>**) Weak Symmetry:  $I_{\cap}(X_1, \dots, X_n : Y)$  is invariant under reordering of  $X_1, \dots, X_n$ .

**Remark:** If (**S<sub>0</sub>**) is satisfied, the first argument of  $I_{\cap}(X_1, \dots, X_n : Y)$  can be treated as a *set* of random variables rather than a *list*. In this case, the notation  $I_{\cap}(\{X_1, \dots, X_n\} : Y)$  would also be appropriate.

For the next set of properties,  $\mathcal{I}(X : Y)$  is a given normative measure of information between  $X$  and  $Y$ . For example,  $\mathcal{I}(X : Y)$  could denote the Shannon mutual information; i.e.,  $\mathcal{I}(X : Y) = I(X : Y)$ . Alternatively, as discussed in the next section, we might take  $\mathcal{I}(X : Y)$  to be the zero-error information. Yet other possibilities for  $\mathcal{I}(X : Y)$  include the Wyner common information [38] or the quantum mutual information [8]. The following are desired properties of intersection information *relative to* the given information measure  $\mathcal{I}$ .

(**LB**) Lowerbound: If  $Q \preceq X_i$  for all  $i = 1, \dots, n$ , then  $I_{\cap}(X_1, \dots, X_n : Y) \geq \mathcal{I}(Q : Y)$ . Under a mild assumption,<sup>1</sup> this equates to  $I_{\cap}(X_1, \dots, X_n : Y) \geq \mathcal{I}(X_1 \wedge \dots \wedge X_n : Y)$ .

(**SR**) Self-Redundancy:  $I_{\cap}(X_1 : Y) = \mathcal{I}(X_1 : Y)$ . The intersection information a single predictor  $X_1$  conveys about the target  $Y$  is equal to the information between the predictor and the target given by the information measure  $\mathcal{I}$ .

(**Id**) Identity:  $I_{\cap}(X, Y : X \vee Y) = \mathcal{I}(X : Y)$ .

(**LP<sub>0</sub>**) Weak Local Positivity:  $I_{\cap}(X_1, X_2 : Y) \geq \mathcal{I}(X_1 : Y) + \mathcal{I}(X_2 : Y) - \mathcal{I}(X_1 \vee X_2 : Y)$ . In other words, for  $n = 2$  predictors, the derived “partial informations” defined in [36] are nonnegative when both (**LP<sub>0</sub>**) and (**GP**) hold.

Finally, we have the less obvious “strong” properties.

(**M<sub>1</sub>**) Strong Monotonicity:  $I_{\cap}(X_1, \dots, X_n, W : Y) \leq I_{\cap}(X_1, \dots, X_n : Y)$  with equality if there exists  $Z \in \{X_1, \dots, X_n, Y\}$  such that  $Z \preceq W$ .

(**S<sub>1</sub>**) Strong Symmetry:  $I_{\cap}(X_1, \dots, X_n : Y)$  is invariant under reordering of  $X_1, \dots, X_n, Y$ .

(**LP<sub>1</sub>**) Strong Local Positivity: For all  $n$ , the derived “partial informations” defined in [36] are non-negative.

---

<sup>1</sup>See Lemmas 3.C.1 and 3.C.2 in Appendix 3.C.1.

Properties **(Eq)**, **(LB)**, and **(M<sub>1</sub>)** are novel and are introduced for the first time here. Given  $I_\cap$ ,  $X_1, \dots, X_n$ ,  $Y$ , and  $Z$ , we define the conditional  $I_\cap$  as:

$$I_\cap(X_1, \dots, X_n : Z | Y) \equiv I_\cap(X_1, \dots, X_n : Y \vee Z) - I_\cap(X_1, \dots, X_n : Y) .$$

This definition of  $I_\cap(X_1, \dots, X_n : Z | Y)$  gives rise to the familiar “chain rule”:

$$I_\cap(X_1, \dots, X_n : Y \vee Z) = I_\cap(X_1, \dots, X_n : Y) + I_\cap(X_1, \dots, X_n : Z | Y) .$$

Some provable<sup>2</sup> properties are:

- $I_\cap(X_1, \dots, X_n : Z | Y) \geq 0$ .
- $I_\cap(X_1, \dots, X_n : Z | Y) = I_\cap(X_1, \dots, X_n : Z)$  if  $Y$  is a constant.

## 3.4 Candidate Intersection Information for Zero-Error Information

### 3.4.1 Zero-Error Information

Introduced in [37], the *zero-error information*, or *Gács-Körner common information*, is a stricter variant of Shannon mutual information. Whereas the mutual information  $I(A : B)$  quantifies the magnitude of information  $A$  conveys about  $B$  with an arbitrarily small error  $\epsilon > 0$ , the zero-error information, denoted  $I^0(A : B)$ , quantifies the magnitude of information  $A$  conveys about  $B$  with *exactly zero* error, i.e.,  $\epsilon = 0$ . The zero-error information between  $A$  and  $B$  equals the entropy of the *common random variable*  $A \wedge B$ ,

$$I^0(A : B) \equiv H(A \wedge B) .$$

An algorithm for computing an instance of the common random variable between two random variables is provided in [37], and straightforwardly generalizes to  $n$  random variables.<sup>3</sup>

Zero-error information has several notable properties, but the most salient is that it is nonnegative and bounded by the mutual information,

$$0 \leq I^0(A : B) \leq I(A : B) .$$

---

<sup>2</sup>See Lemma 3.C.3 in Appendix 3.C.1.

<sup>3</sup>See Appendix 3.A.

This generalizes to arbitrary  $n$ :

$$0 \leq I^0(X_1 : \dots : X_n) \leq \min_{i,j} I(X_i : X_j) .$$

### 3.4.2 Intersection Information for Zero-Error Information

It is pleasingly straightforward to define a palatable intersection information for zero-error information (i.e., setting  $\mathcal{I} = I^0$  as the normative measure of information). We propose the zero-error intersection information,  $I_\lambda^0(X_1, \dots, X_n : Y)$ , as the maximum zero-error information  $I^0(Q : Y)$  that some random variable  $Q$  conveys about  $Y$ , subject to  $Q$  being a function of each predictor  $X_1, \dots, X_n$ :

$$\begin{aligned} I_\lambda^0(X_1, \dots, X_n : Y) &\equiv \max_{\Pr(Q|Y)} I^0(Q : Y) \\ &\text{subject to } \forall i \in \{1, \dots, n\} : Q \preceq X_i . \end{aligned} \quad (3.2)$$

Basic algebra<sup>4</sup> shows that a maximizing  $Q$  is the common random variable across all predictors. This substantially simplifies eq. (3.2) to:

$$\begin{aligned} I_\lambda^0(X_1, \dots, X_n : Y) &= I^0(X_1 \wedge \dots \wedge X_n : Y) \\ &= H[(X_1 \wedge \dots \wedge X_n) \wedge Y] \\ &= H(X_1 \wedge \dots \wedge X_n \wedge Y) . \end{aligned} \quad (3.3)$$

Importantly, the zero-error information,  $I_\lambda^0(X_1, \dots, X_n : Y)$  satisfies ten of the twelve desired properties from Section 3.3.4, leaving only  $(\mathbf{LP}_0)$  and  $(\mathbf{LP}_1)$  unsatisfied.<sup>5</sup>

## 3.5 Candidate Intersection Information for Shannon Information

In the last section, we defined an intersection information for zero-error information which satisfies the vast majority of desired properties. This is a solid start, but an intersection information for Shannon mutual information remains the goal. Towards this end, we use the same method as in eq. (3.2), leading to  $I_\lambda$ , our candidate intersection information measure for Shannon mutual information,

$$\begin{aligned} I_\lambda(X_1, \dots, X_n : Y) &\equiv \max_{\Pr(Q|Y)} I(Q : Y) \\ &\text{subject to } Q \preceq X_i \ \forall i \in \{1, \dots, n\} . \end{aligned} \quad (3.4)$$

---

<sup>4</sup>See Lemma 3.D.1 in Appendix 5.D.

<sup>5</sup>See Lemmas 3.C.4, 3.C.5, 3.C.6 in Appendix 3.C.2.



With some algebra<sup>6</sup> this similarly simplifies to,

$$I_{\lambda}(X_1, \dots, X_n : Y) = I(X_1 \wedge \dots \wedge X_n : Y) . \quad (3.5)$$

Unfortunately  $I_{\lambda}$  does not satisfy as many of the desired properties as  $I_{\lambda}^0$ . However, our candidate  $I_{\lambda}$  still satisfies 7 of the 12 properties (Table 3.1), most importantly the enviable **(TM)**,<sup>7</sup> which has, until now, not been satisfied by any proposed measure. Table 3.1 lists the desired properties satisfied by  $I_{\min}$ ,  $I_{\lambda}$ , and  $I_{\lambda}^0$ . For reference, we also include  $I_{\text{red}}$ , the proposed measure from [18].

Comparing the three subject intersection information measures,<sup>8</sup> we have:

$$0 \leq I_{\lambda}^0(X_1, \dots, X_n : Y) \leq I_{\lambda}(X_1, \dots, X_n : Y) \leq I_{\min}(X_1, \dots, X_n : Y) . \quad (3.6)$$

	Property	$I_{\min}$	$I_{\text{red}}$	$I_{\lambda}$	$I_{\lambda}^0$
<b>(GP)</b>	Global Positivity	✓	✓	✓	✓
<b>(Eq)</b>	Equivalence-Class Invariance	✓	✓	✓	✓
<b>(TM)</b>	Target Monotonicity			✓	✓
<b>(M<sub>0</sub>)</b>	Weak Monotonicity	✓		✓	✓
<b>(S<sub>0</sub>)</b>	Weak Symmetry	✓	✓	✓	✓
<b>(LB)</b>	Lowerbound	✓	✓	✓	✓
<b>(SR)</b>	Self-Redundancy	✓	✓	✓	✓
<b>(Id)</b>	Identity		✓		✓
<b>(LP<sub>0</sub>)</b>	Weak Local Positivity	✓	✓		
<b>(M<sub>1</sub>)</b>	Strong Monotonicity				✓
<b>(S<sub>1</sub>)</b>	Strong Symmetry				✓
<b>(LP<sub>1</sub>)</b>	Strong Local Positivity	✓			

Table 3.1: The  $I_{\cap}$  desired properties each measure satisfies.

Despite not satisfying **(LP<sub>0</sub>)**,  $I_{\lambda}$  remains an important stepping-stone towards the ideal Shannon  $I_{\cap}$ . First,  $I_{\lambda}$  captures what is inarguably redundant information (the common random variable); this makes  $I_{\lambda}$  necessarily a lower bound on any reasonable redundancy measure. Second, it is the first proposal to satisfy target monotonicity and the associated chain rule. Lastly,  $I_{\lambda}$  is the first measure to reach intuitive answers in many canonical situations, while also being generalizable to an arbitrary number of inputs.

<sup>6</sup>See Lemma 3.D.2 in Appendix 5.D.

<sup>7</sup>See Lemmas 3.C.7, 3.C.8, 3.C.9 in Appendix 3.C.3.

<sup>8</sup>See Lemma 3.D.3 in Appendix 5.D.

### 3.6 Three Examples Comparing $I_{\min}$ and $I_{\lambda}$

Examples UNQ and RDNXOR illustrate  $I_{\lambda}$ 's successes, and example IMPERFECTRDN illustrates  $I_{\lambda}$ 's paramount deficiency. For each example we show the joint distribution  $\Pr(x_1, x_2, y)$ , a diagram, and the decomposition derived from setting  $I_{\min}/I_{\lambda}$  as the  $I_{\cap}$  measure. At each lattice junction, the left number is the  $I_{\cap}$  value of that node, and the number in parentheses is the  $I_{\partial}$  value.<sup>9</sup> Readers unfamiliar with the  $n = 2$  partial information lattice should consult [36], but in short,  $I_{\partial}$  measures the amount of “new” information at this node in the lattice compared to nodes lower in the lattice. Except for IMPERFECTRDN, measures  $I_{\lambda}$  and  $I_{\lambda}^0$  reach the same decomposition for all presented examples. Per [36], the four partial informations are calculated as follows:

$$\begin{aligned}
 I_{\partial}(X_1, X_2 : Y) &= I_{\cap}(X_1, X_2 : Y) \\
 I_{\partial}(X_1 : Y) &= I(X_1 : Y) - I_{\cap}(X_1, X_2 : Y) \\
 I_{\partial}(X_2 : Y) &= I(X_2 : Y) - I_{\cap}(X_1, X_2 : Y) \\
 I_{\partial}(X_1 \vee X_2 : Y) &= I(X_1 \vee X_2 : Y) - I(X_1 : Y) - I(X_2 : Y) + I_{\cap}(X_1, X_2 : Y) \\
 &= I(X_1 \vee X_2 : Y) - I_{\partial}(X_1 : Y) - I_{\partial}(X_2 : Y) - I_{\partial}(X_1, X_2 : Y) .
 \end{aligned} \tag{3.7}$$

**Example Unq** (Figure 3.3). The desired decomposition for this example is two bits of unique information;  $X_1$  uniquely specifies one bit of  $Y$ , and  $X_2$  uniquely specifies the other bit of  $Y$ . The chief criticism of  $I_{\min}$  in [15] was that  $I_{\min}$  calculated one bit of redundancy and one bit of synergy for UNQ (Figure 3.3c). We see that unlike  $I_{\min}$ ,  $I_{\lambda}$  satisfyingly arrives at two bits of unique information. This is easily seen by the inequality,

$$0 \leq I_{\lambda}(X_1, X_2 : Y) \leq H(X_1 \wedge X_2) \leq I(X_1 : X_2) = 0 \text{ bits} . \tag{3.8}$$

Therefore, as  $I(X_1 : X_2) = 0$ , we have  $I_{\lambda}(X_1, X_2 : Y) = 0$  bits leading to  $I_{\partial}(X_1 : Y) = 1$  bit and  $I_{\partial}(X_2 : Y) = 1$  bit (Figure 3.3d).

**Example RdnXor** (Figure 3.4). In [15], RDNXOR was an example where  $I_{\min}$  shined by reaching the desired decomposition of one bit of redundancy and one bit of synergy. We see that  $I_{\lambda}$  finds this same answer.  $I_{\lambda}$  extracts the common random variable within  $X_1$  and  $X_2$ , the  $\mathbf{r}/\mathbf{R}$  bit, and calculates the mutual information between the common random variable and  $Y$  to arrive at  $I_{\lambda}(X_1, X_2 : Y) = 1$  bit.

**Example ImperfectRdn** (Figure 3.5). IMPERFECTRDN highlights the foremost shortcoming of  $I_{\lambda}$ ;  $I_{\lambda}$  does not detect “imperfect” or “lossy” correlations between  $X_1$  and  $X_2$ . Given  $(\mathbf{LP}_0)$ ,

---

<sup>9</sup>This is the same notation used in [7].

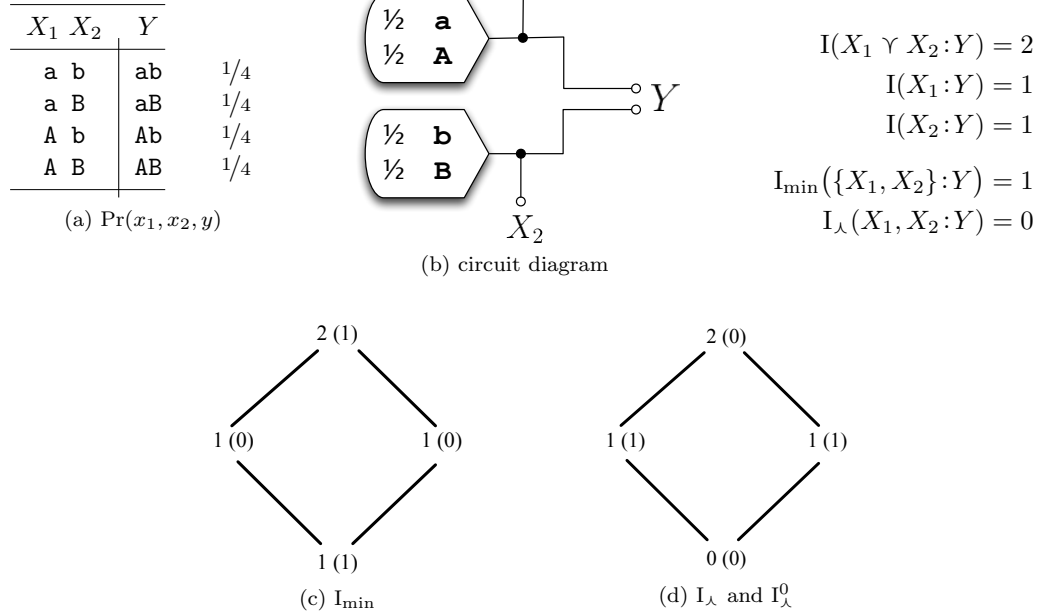


Figure 3.3: Example UNQ. This is the canonical example of unique information.  $X_1$  and  $X_2$  each uniquely specify a single bit of  $Y$ . This is the simplest example where  $I_{\min}$  calculates an undesirable decomposition (c) of one bit of redundancy and one bit of synergy.  $I_{\lambda}$  and  $I_{\lambda}^0$  each calculate the desired decomposition (d).

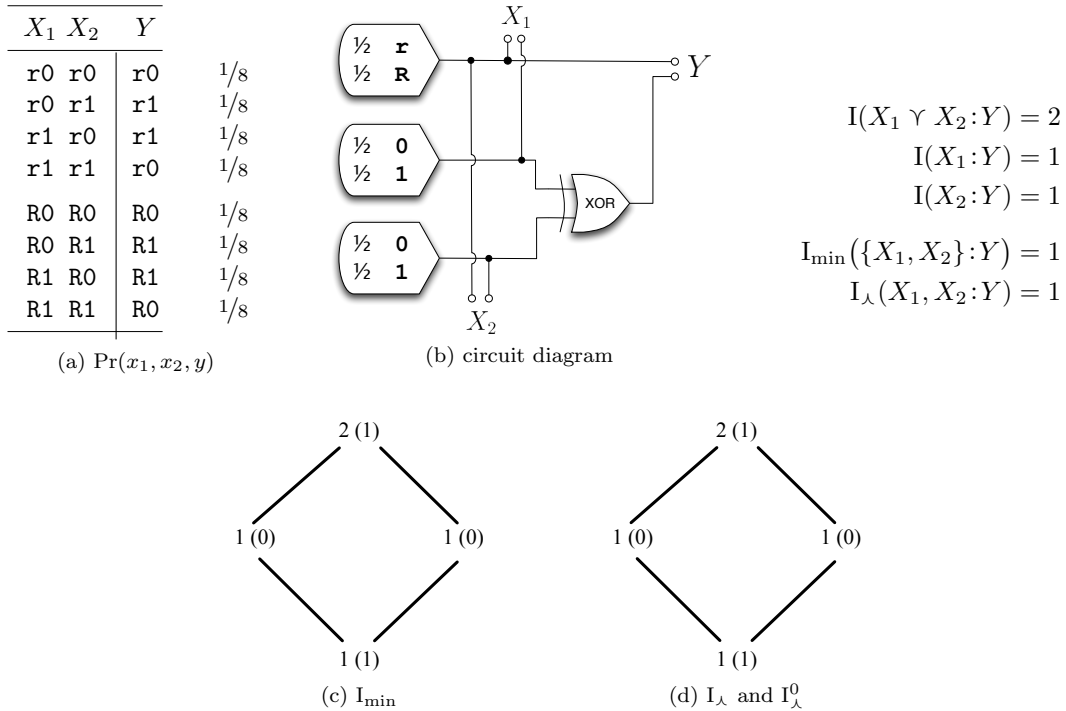


Figure 3.4: Example RDNXOR. This is the canonical example of redundancy and synergy coexisting.  $I_{\min}$  and  $I_{\lambda}$  each reach the desired decomposition of one bit of redundancy and one bit of synergy. This is the simplest example demonstrating  $I_{\lambda}$  and  $I_{\lambda}^0$  correctly extracting the embedded redundant bit within  $X_1$  and  $X_2$ .

we can determine the desired decomposition analytically. First,  $I(X_1 \vee X_2 : Y) = I(X_1 : Y) = 1$  bit; therefore,  $I(X_2 : Y | X_1) = I(X_1 \vee X_2 : Y) - I(X_1 : Y) = 0$  bits. This determines two of the partial informations—the synergistic information  $I_\partial(X_1 \vee X_2 : Y)$  and the unique information  $I_\partial(X_2 : Y)$  are both zero. Then, the redundant information  $I_\partial(X_1, X_2 : Y) = I(X_2 : Y) - I_\partial(X_2 : Y) = I(X_2 : Y) = 0.99$  bits. Having determined three of the partial informations, we compute the final unique information  $I_\partial(X_1 : Y) = I(X_1 : Y) - 0.99 = 0.01$  bits.

How well do  $I_{\min}$  and  $I_\lambda$  match the desired decomposition of IMPERFECTRDN? We see that  $I_{\min}$  calculates the desired decomposition (Figure 3.5c); however,  $I_\lambda$  does not (Figure 3.5d). Instead,  $I_\lambda$  calculates zero redundant information, that  $I_\cap(X_1, X_2 : Y) = 0$  bits. This unpleasant answer arises from  $\Pr(X_1 = 1, X_2 = 0) > 0$ . If this were zero, IMPERFECTRDN reverts to the example RDN (Figure ?? in Appendix 3.E) where both  $I_\lambda$  and  $I_{\min}$  reach the desired one bit of redundant information. Due to the nature of the common random variable,  $I_\lambda$  only sees the “deterministic” correlations between  $X_1$  and  $X_2$ —add even an iota of noise between  $X_1$  and  $X_2$  and  $I_\lambda$  plummets to zero. This highlights a related issue with  $I_\lambda$ —it is not continuous; an arbitrarily small change in the probability distribution can result in a discontinuous jump in the value of  $I_\lambda$ . As with traditional information measures, such as the entropy and the mutual information, it may be desirable to have an  $I_\cap$  measure that is continuous over the simplex.

To summarize, IMPERFECTRDN shows that when there are additional “imperfect” correlations between  $A$  and  $B$ , i.e.  $I(A : B | A \wedge B) > 0$ ,  $I_\lambda$  sometimes *underestimates* the ideal  $I_\cap(A, B : Y)$ .

### 3.7 Negative synergy and state-dependent (GP)

In IMPERFECTRDN we saw  $I_\lambda$  calculate a synergy of  $-0.99$  bits (Figure 3.5d). What does this mean? Could negative synergy be a “real” property of Shannon information? When  $n = 2$ , it’s fairly easy to diagnose the cause of negative synergy from the equation for  $I_\partial(X_1, X_2 : Y)$  in eq. (3.7). Given (GP) and (SR), negative synergy occurs if and only if,

$$\begin{aligned} I(X_1 \vee X_2 : Y) &< I(X_1 : Y) + I(X_2 : Y) - I_\cap(X_1, X_2 : Y) \\ &= I_\cup(X_1, X_2 : Y) . \end{aligned} \tag{3.9}$$

From eq. (3.9), we see negative synergy occurs when  $I_\cap$  is small, perhaps *too small*. Equivalently, negative synergy occurs when the joint r.v. conveys *less* about  $Y$  than the two r.v.’s  $X_1$  and  $X_2$  convey separately—mathematically, when  $I(X_1 \vee X_2 : Y) < I_\cup(X_1, X_2 : Y)$ .<sup>10</sup> On the face of it this sounds strange. No usable structure in  $X_1$  or  $X_2$  “disappears” after they are combined by

<sup>10</sup> $I_\cap$  and  $I_\cup$  are duals related by the inclusion–exclusion principle. For arbitrary  $n$ , this is  $I_\cup(X_1, \dots, X_n : Y) = \sum_{\mathbf{S} \subseteq \{X_1, \dots, X_n\}} (-1)^{|\mathbf{S}|+1} I_\cap(S_1, \dots, S_{|\mathbf{S}|} : Y)$ .

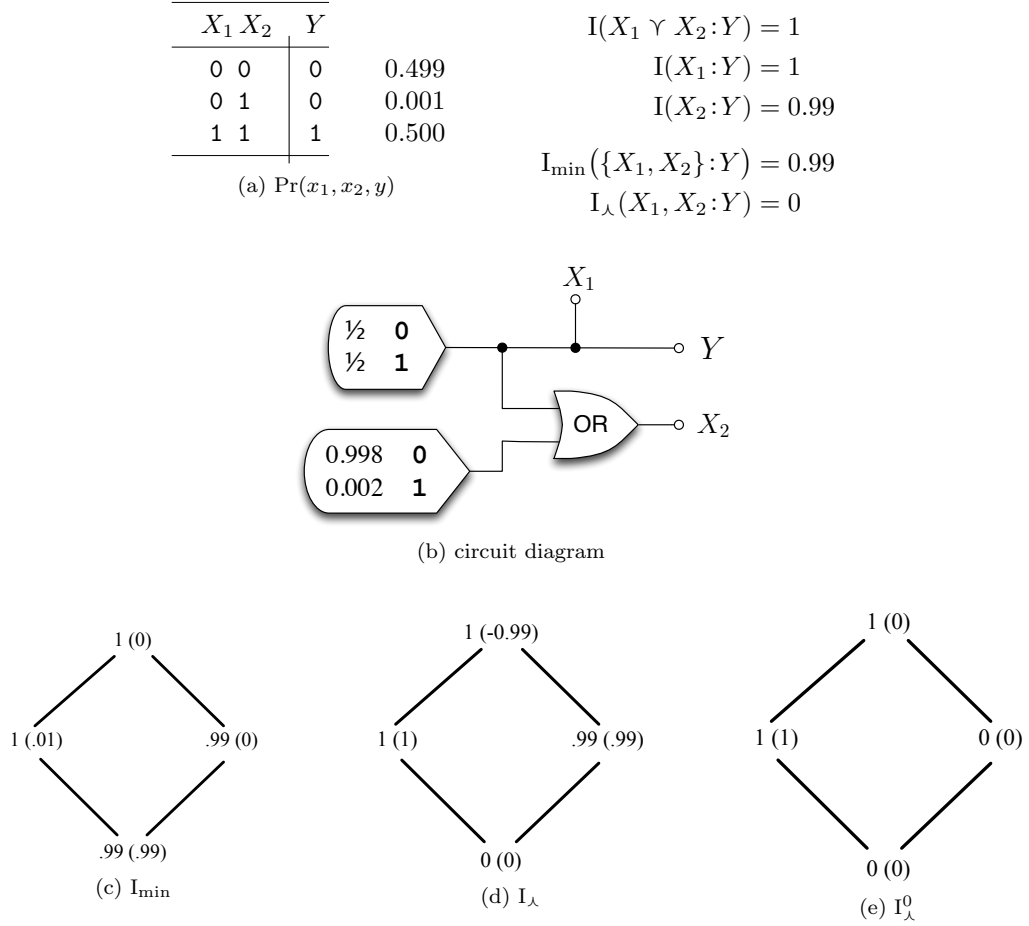


Figure 3.5: Example IMPERFECTRDN.  $I_{\lambda}$  is blind to the noisy correlation between  $X_1$  and  $X_2$  and calculates zero redundant information. An ideal  $I_{\cap}$  measure would detect that all of the information  $X_2$  specifies about  $Y$  is also specified by  $X_1$  to calculate  $I_{\cap}(X_1, X_2 : Y) = 0.99$  bits.

$Z = X_1 \vee X_2$ . By the definition of  $\vee$ , there are always functions  $f_1$  and  $f_2$  such that  $X_1 \cong f_1(Z)$  and  $X_2 \cong f_2(Z)$ . Therefore, if your favorite  $I_\cap$  measure does not satisfy **(LP<sub>0</sub>)**, it is likely too strict.

This means that, to our surprise, our measure  $I_\lambda^0$  does not account for the full zero-information overlap between  $I^0(X_1:Y)$  and  $I^0(X_2:Y)$ . This is shown in example SUBTLE (Figure 3.6) where  $I_\lambda^0$  calculates a synergy of  $-0.252$  bits. Defining a zero-error  $I_\cap$  that satisfies **(LP<sub>0</sub>)** is a matter of ongoing research.

### 3.7.1 Consequences of state-dependent **(GP)**

In [15] it's argued that  $I_{\min}$  upperbounds the ideal  $I_\cap$ . Inspired by  $I_{\min}$  assuming state-dependent **(SR)** and **(M<sub>0</sub>)** to achieve a tighter upperbound on  $I_\cap$ , we assume state-dependent **(GP)** to achieve a tighter lowerbound on  $I_\cap$  for  $n = 2$ . Our bound, denoted  $I_{\text{smp}}$  for “sum minus pair”, is defined as,

$$I_{\text{smp}}(X_1, X_2 : Y) \equiv \sum_{y \in Y} \Pr(y) \max [0, I(X_1:y) + I(X_2:y) - I(X_1 \vee X_2:y)] , \quad (3.10)$$

where  $I(\bullet:y)$  is the same Kullback-Liebler divergence from eq. (3.1).

For example SUBTLE, the target  $Y \cong X_1 \vee X_2$ , therefore per **(Id)**,  $I_\cap(X_1, X_2:Y) = I(X_1:X_2) = 0.252$  bits. However, given state-dependent **(GP)**, applying  $I_{\text{smp}}$  yields  $I_\cap(X_1, X_2:Y) \geq 0.390$ . Therefore, **(Id)** and state-dependent **(GP)** are incompatible. Secondly, given state-dependent **(GP)**, example SUBTLE additionally illustrates a conjecture from [7] that the intersection information two predictors have about a target can exceed the mutual information between them, i.e.,  $I_\cap(X_1, X_2:Y) \not\leq I(X_1:X_2)$ .

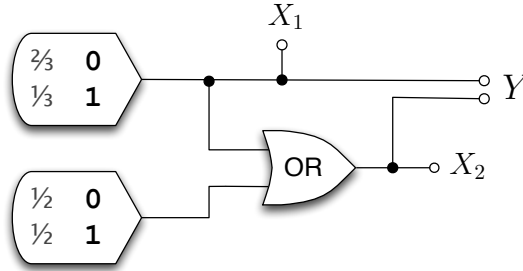
## 3.8 Conclusion and Path Forward

We've made incremental progress on several fronts towards the ideal Shannon  $I_\cap$ .

**Desired Properties.** We have tightened, expanded, and pruned the desired properties for  $I_\cap$ . Particularly,

- **(LB)** is a non-contentious yet tighter lower-bound on  $I_\cap$  than **(GP)**.
- Motivated by the natural equality  $I_\cap(X_1, \dots, X_n:Y) = I_\cap(X_1, \dots, X_n, Y:Y)$ , we introduce **(M<sub>1</sub>)** as a desired property.
- What was before an implicit assumption, we introduce **(Eq)** to better ground one's thinking.
- A separate chain-rule property is superfluous. Any desirable properties of conditional  $I_\cap$  are simply consequences of **(GP)** and **(TM)**.

$X_1$	$X_2$	$Y$		
0	0	00	$1/3$	$I(X_1 \vee X_2 : Y) = 1.585$
0	1	01	$1/3$	$I(X_1 : Y) = 0.918$
1	1	11	$1/3$	$I(X_2 : Y) = 0.918$
				$I(X_1 : X_2) = 0.252$
(a) $\Pr(x_1, x_2, y)$				$I_{\min}(\{X_1, X_2\} : Y) = 0.585$
				$I_{\lambda}(X_1, X_2 : Y) = 0.0$
				$I_{\text{smp}}(X_1, X_2 : Y) = 0.390$



(b) circuit diagram

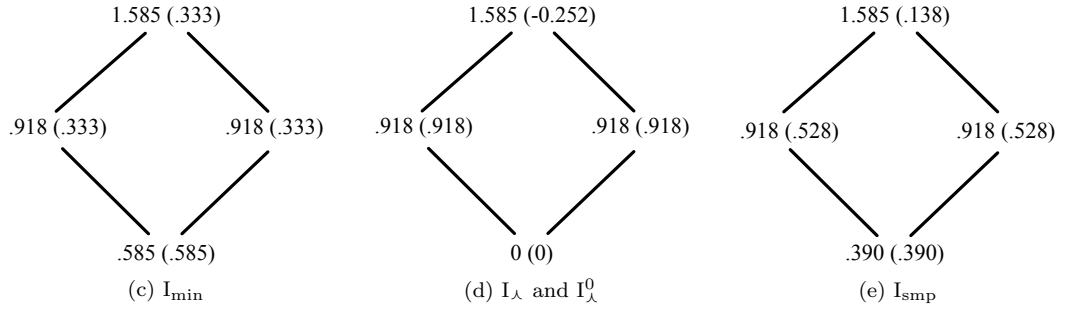


Figure 3.6: Example SUBTLE. In this example both  $I_{\lambda}$  and  $I_{\lambda}^0$  calculate a synergy of  $-0.252$  bits of synergy. What kind of redundancy must be captured for a nonnegative decomposition for this example?

**A new measure.** Based on the Gács-Körner common random variable, we introduced a new Shannon  $I_{\cap}$  measure. Our measure,  $I_{\lambda}$ , is theoretically principled and the first to satisfy **(TM)**.

**How to improve.** We identified where  $I_{\lambda}$  fails; it does not detect “imperfect” correlations between  $X_1$  and  $X_2$ . One next step is to develop a less stringent  $I_{\cap}$  measure that satisfies **(LP<sub>0</sub>)** for simple nondeterministic examples like IMPERFECTRDN while still satisfying **(TM)**.

To our surprise, example SUBTLE shows that  $I_{\lambda}^0$  does not satisfy **(LP<sub>0</sub>)**! This suggests that  $I_{\lambda}^0$  is too strict—what kind of zero-error informational overlap is  $I_{\lambda}^0$  not capturing? A separate next step is to formalize what exactly is required for a zero-error  $I_{\cap}$  to satisfy **(LP<sub>0</sub>)**. From SUBTLE we can likewise see that for zero-error information, **(LP<sub>0</sub>)** is incompatible with **(Id)**.

Finally, we showed that state-dependent **(GP)**, a seemingly reasonable property, is incompatible with **(Id)** and moreover entails that  $I_{\cap}(X_1, X_2 : Y)$  can exceed  $I(X_1 : X_2)$ .



# Appendix

## 3.A Algorithm for Computing Common Random Variable

Given  $n$  random variables  $X_1, \dots, X_n$ , the common random variable  $X_1 \wedge \dots \wedge X_n$  is computed by steps 1–3 in Appendix 3.B.

## 3.B Algorithm for Computing $I_\wedge$

1. For each  $X_i$  for  $i = 1, \dots, n$ , take its states  $x_i$  and place them as nodes on a graph. At the end of this process there will be  $\sum_{i=1}^n |X_i|$  nodes on the graph.
2. For each pair of RVs  $X_i, X_j$  ( $i \neq j$ ), draw an undirected edge connecting nodes  $x_i$  and  $x_j$  if  $\Pr(x_i, x_j) > 0$ . At the end of this process the undirected graph will consist of  $k$  connected components  $1 \leq k \leq \min_i |X_i|$ . Denote these  $k$  disjoint components as  $\mathbf{c}_1, \dots, \mathbf{c}_k$ .
3. Each connected component of the graph constitutes a distinct state of the common random variable  $Q$ , i.e.,  $|Q| = k$ . Denote the states of the common random variable  $Q$  by  $q_1, \dots, q_k$ .
4. Construct the joint probability distribution  $\Pr(Q, Y)$  as follows. For every state  $(q_i, y) \in Q \times Y$ , the joint probability is created by summing over the entries of  $\Pr(x_1, \dots, x_n, y)$  in component  $i$ . More precisely,

$$\Pr(Q = q_i, Y = y) = \sum_{x_1, \dots, x_n} \Pr(x_1, \dots, x_n, y) \quad \text{if } \{x_1, \dots, x_n\} \subseteq \mathbf{c}_i .$$

5. Using  $\Pr(Q, Y)$ , compute  $I_\wedge(X_1, \dots, X_n : Y)$  simply by computing the Shannon mutual information between  $Q$  and  $Y$ , i.e.,  $I(Q : Y) = D_{\text{KL}} [\Pr(Q, Y) \| \Pr(Q) \Pr(Y)]$ .

## 3.C Lemmas and Proofs

### 3.C.1 Lemmas on Desired Properties

**Lemma 3.C.1.** *If (LB) holds, then  $I_\cap(X_1, \dots, X_n : Y) \geq \mathcal{I}(X_1 \wedge \dots \wedge X_n : Y)$ .*

*Proof.* Assume that **(LB)** holds. By definition,  $X_1 \wedge \cdots \wedge X_n \preceq X_i$  for  $i = 1, \dots, n$ . So, by **(LB)**, we immediately conclude that  $I_\cap(X_1, \dots, X_n : Y) \geq \mathcal{I}(X_1 \wedge \cdots \wedge X_n : Y)$ , which is the desired result.  $\square$

For the converse, we need the following assumption:

**(IM)** If  $X_1 \preceq X_2$ , then  $\mathcal{I}(X_1 : Y) \leq \mathcal{I}(X_2 : Y)$ .

**Lemma 3.C.2.** *Suppose that **(IM)** holds, and that  $I_\cap(X_1, \dots, X_n : Y) \geq \mathcal{I}(X_1 \wedge \cdots \wedge X_n : Y)$ . Then **(LB)** holds.*

*Proof.* Assume that  $I_\cap(X_1, \dots, X_n : Y) \geq \mathcal{I}(X_1 \wedge \cdots \wedge X_n : Y)$ . Let  $Q \preceq X_i$  for  $i = 1, \dots, n$ . Because  $X_1 \wedge \cdots \wedge X_n$  is the largest (informationally richest) random variable that is informationally poorer than  $X_i$  for  $i = 1, \dots, n$ , it follows that  $Q \preceq X_1 \wedge \cdots \wedge X_n$ . Therefore, by **(IM)**,  $\mathcal{I}(X_1 \wedge \cdots \wedge X_n : Y) \geq \mathcal{I}(Q : Y)$ . Hence,  $I_\cap(X_1, \dots, X_n : Y) \geq \mathcal{I}(Q : Y)$  also, which completes the proof.  $\square$

**Remark:** Assumption **(IM)** is satisfied by zero-error information and Shannon mutual information.

**Lemma 3.C.3.** *Given  $I_\cap$ ,  $X_1, \dots, X_n$ ,  $Y$ , and  $Z$ , consider the conditional intersection information*

$$I_\cap(X_1, \dots, X_n : Z|Y) = I_\cap(X_1, \dots, X_n : Y \vee Z) - I_\cap(X_1, \dots, X_n : Y).$$

*Suppose that **(GP)**, **(Eq)**, and **(TM)** hold. Then, the following properties hold:*

- $I_\cap(X_1, \dots, X_n : Z|Y) \geq 0$ .
- $I_\cap(X_1, \dots, X_n : Z|Y) = I_\cap(X_1, \dots, X_n : Z)$  if  $Y$  is a constant.

*Proof.* We have  $Y \preceq Y \vee Z$ . Therefore, by **(TM)**, it immediately follows that  $I_\cap(X_1, \dots, X_n : Z|Y) \geq 0$ .

Next, suppose that  $Y$  is a constant. Then  $Y \preceq Z$ , and hence  $Y \vee Z \cong Z$ . By **(Eq)**,  $I_\cap(X_1, \dots, X_n : Y \vee Z) = I_\cap(X_1, \dots, X_n : Z)$ . Moreover, by **(GP)**,  $I_\cap(X_1, \dots, X_n : Y) = 0$ . Thus,  $I_\cap(X_1, \dots, X_n : Z|Y) = I_\cap(X_1, \dots, X_n : Z)$  as desired.  $\square$

### 3.C.2 Properties of $I_\lambda^0$

**Lemma 3.C.4.** *The measure of intersection information  $I_\lambda^0(X_1, \dots, X_n : Y)$  satisfies **(GP)**, **(Eq)**, **(TM)**, **(M<sub>0</sub>)**, and **(S<sub>0</sub>)**, but not **(LP<sub>0</sub>)**.*

*Proof.* (**GP**): The inequality  $I_\lambda^0(X_1, \dots, X_n : Y) \geq 0$  follows immediately from the identity  $I_\lambda^0(X_1, \dots, X_n : Y) = H(X_1 \wedge \dots \wedge X_n \wedge Y)$  and the nonnegativity of  $H(\cdot)$ . Next, if  $Y$  is a constant, then by the generalized absorption law,  $X_1 \wedge \dots \wedge X_n \wedge Y \cong Y$ . Thus, by the invariance of  $H(\cdot)$  (Lemma 3.3.1(a)),  $H(X_1 \wedge \dots \wedge X_n \wedge Y) = H(Y) = 0$ .

(**Eq**): Consider  $X_1 \wedge \dots \wedge X_n \wedge Y$ . The equivalence class (with respect to  $\cong$ ) in which this random variable resides is closed under substitution of  $X_i$  (for  $i = 1, \dots, n$ ) or  $Y$  by an informationally equivalent random variable. Hence, because  $I_\lambda^0(X_1, \dots, X_n : Y) = H(X_1 \wedge \dots \wedge X_n \wedge Y)$  and  $H(\cdot)$  is invariant over the equivalence class of random variables that are informationally equivalent to  $X_1 \wedge \dots \wedge X_n \wedge Y$  (by Lemma 3.3.1(a)), the desired result holds.

(**TM**): Suppose that  $Y \preceq Z$ . Then,  $X_1 \wedge \dots \wedge X_n \wedge Y \preceq X_1 \wedge \dots \wedge X_n \wedge Z$ . Then, we have

$$\begin{aligned} I_\lambda^0(X_1, \dots, X_n : Y) &= H(X_1 \wedge \dots \wedge X_n \wedge Y) \\ &\leq H(X_1 \wedge \dots \wedge X_n \wedge Z) \quad \text{by monotonicity of } H(\cdot) \text{ (Lemma 3.3.1(b))} \\ &= I_\lambda^0(X_1, \dots, X_n : Z), \end{aligned}$$

as desired.

(**M0**): By the generalized absorption law,  $X_1 \wedge \dots \wedge X_n \wedge W \wedge Y \preceq X_1 \wedge \dots \wedge X_n \wedge Y$ . Hence,

$$\begin{aligned} I_\lambda^0(X_1, \dots, X_n, W : Y) &= H(X_1 \wedge \dots \wedge X_n \wedge W \wedge Y) \\ &\leq H(X_1 \wedge \dots \wedge X_n \wedge Y) \quad \text{by monotonicity of } H(\cdot) \text{ (Lemma 3.3.1(b))} \\ &= I_\lambda^0(X_1, \dots, X_n : Y), \end{aligned}$$

as desired.

Next, suppose that there exists  $Z \in \{X_1, \dots, X_n\}$  such that  $Z \preceq W$ . Then, by the generalized absorption law,  $X_1 \wedge \dots \wedge X_n \wedge W \wedge Y \cong X_1 \wedge \dots \wedge X_n \wedge Y$ . Hence,

$$\begin{aligned} I_\lambda^0(X_1, \dots, X_n, W : Y) &= H(X_1 \wedge \dots \wedge X_n \wedge W \wedge Y) \\ &= H(X_1 \wedge \dots \wedge X_n \wedge Y) \quad \text{by invariance of } H(\cdot) \text{ (Lemma 3.3.1(a))} \\ &= I_\lambda^0(X_1, \dots, X_n : Y), \end{aligned}$$

as desired.

(**S0**): By the commutativity law,  $X_1 \wedge \dots \wedge X_n \wedge Y$  is invariant (with respect to  $\cong$ ) under reordering of  $X_1, \dots, X_n$ . Hence, the desired result follows immediately from the identity  $I_\lambda^0(X_1, \dots, X_n : Y) = H(X_1 \wedge \dots \wedge X_n \wedge Y)$  and the invariance of  $H(\cdot)$  (Lemma 3.3.1(a)).

(**LP**<sub>0</sub>): For  $I_\lambda^0$ , (**LP**<sub>0</sub>) relative to zero-error information can be written as

$$H(X_1 \wedge X_2 \wedge Y) \geq H(X_1 \wedge Y) + H(X_2 \wedge Y) - H((X_1 \vee X_2) \wedge Y). \quad (3.11)$$

However, this inequality does not hold in general. To see this, suppose that it does hold for arbitrary  $X_1$ ,  $X_2$ , and  $Y$ . Note that  $(X_1 \vee X_2) \wedge Y \preceq Y$ , which implies that  $H((X_1 \vee X_2) \wedge Y) \leq H(Y)$  (by monotonicity of  $H(\cdot)$ ). Hence, the inequality (3.11) implies that

$$H(X_1 \wedge X_2 \wedge Y) \geq H(X_1 \wedge Y) + H(X_2 \wedge Y) - H(Y).$$

Rewriting this, we get

$$H(X_1 \wedge Y) + H(Y \wedge X_2) \leq H(X_1 \wedge Y \wedge X_2) + H(Y).$$

But this is the supermodularity law for common information, which is known to be false in general; see [22], Section 5.4. □

**Lemma 3.C.5.** *With respect to zero-error information, the measure of intersection information  $I_\lambda^0(X_1, \dots, X_n : Y)$  satisfies (**LB**), (**SR**), and (**Id**).*

*Proof.* (**LB**): Suppose that  $Q \preceq X_i$  for  $i = 1, \dots, n$ . Because  $X_1 \wedge \dots \wedge X_n$  is the largest (informationally richest) random variable that is informationally poorer than  $X_i$  for  $i = 1, \dots, n$ , it follows that  $Q \preceq X_1 \wedge \dots \wedge X_n$ . This implies that  $X_1 \wedge \dots \wedge X_n \wedge Y \succeq Q \wedge Y$ . Therefore,

$$\begin{aligned} I_\lambda^0(X_1, \dots, X_n : Y) &= H(X_1 \wedge \dots \wedge X_n \wedge Y) \\ &\geq H(Q \wedge Y) \quad \text{by monotonicity of } H(\cdot) \text{ (Lemma 3.3.1(b))} \\ &= I^0(Q : Y), \end{aligned}$$

as desired.

(**SR**): We have  $I_\lambda^0(X_1 : Y) = H(X_1 \wedge Y) = I^0(X_1 : Y)$ .

(**Id**): By the associative and absorption laws, we have  $X \wedge Y \wedge (X \vee Y) \cong X \wedge Y$ . Thus,

$$\begin{aligned} I_\lambda^0(X, Y : X \vee Y) &= H(X \wedge Y \wedge (X \vee Y)) \\ &= H(X \wedge Y) \quad \text{by invariance of } H(\cdot) \text{ (Lemma 3.3.1(a))} \\ &= I^0(X : Y), \end{aligned}$$

as desired.

□

**Lemma 3.C.6.** *The measure of intersection information  $I_\lambda^0(X_1, \dots, X_n : Y)$  satisfies  $(\mathbf{M}_1)$  and  $(\mathbf{S}_1)$ , but not  $(\mathbf{LP}_1)$ .*

*Proof.*  $(\mathbf{M}_1)$ : The desired inequality is identical to  $(\mathbf{M}_0)$ , so it remains to prove the sufficient condition for equality. Suppose that there exists  $Z \in \{X_1, \dots, X_n, Y\}$  such that  $Z \preceq W$ . Then, by the generalized absorption law,  $X_1 \wedge \dots \wedge X_n \wedge W \wedge Y \cong X_1 \wedge \dots \wedge X_n \wedge Z$ . Hence,

$$\begin{aligned} I_\lambda^0(X_1, \dots, X_n, W : Y) &= H(X_1 \wedge \dots \wedge X_n \wedge W \wedge Y) \\ &= H(X_1 \wedge \dots \wedge X_n \wedge Z) \quad \text{by invariance of } H(\cdot) \text{ (Lemma 3.3.1(a))} \\ &= I_\lambda^0(X_1, \dots, X_n : Z), \end{aligned}$$

as desired.

$(\mathbf{S}_1)$ : By the commutativity law,  $X_1 \wedge \dots \wedge X_n \wedge Y$  is invariant (with respect to  $\cong$ ) under reordering of  $X_1, \dots, X_n, Y$ . Hence, the desired result follows immediately from the identity  $I_\lambda^0(X_1, \dots, X_n : Y) = H(X_1 \wedge \dots \wedge X_n \wedge Y)$  and the invariance of  $H(\cdot)$  (Lemma 3.3.1(a)).

$(\mathbf{LP}_1)$ : This follows from not satisfying  $(\mathbf{LP}_0)$ .

□

### 3.C.3 Properties of $I_\lambda$

**Lemma 3.C.7.** *The measure of intersection information  $I_\lambda(X_1, \dots, X_n : Y)$  satisfies  $(\mathbf{GP})$ ,  $(\mathbf{Eq})$ ,  $(\mathbf{TM})$ ,  $(\mathbf{M}_0)$ , and  $(\mathbf{S}_0)$ , but not  $(\mathbf{LP}_0)$ .*

*Proof.*  $(\mathbf{GP})$ : The inequality  $I_\lambda(X_1, \dots, X_n : Y) \geq 0$  follows immediately from the identity  $I_\lambda(X_1, \dots, X_n : Y) = I(X_1 \wedge \dots \wedge X_n : Y)$  and the nonnegativity of mutual information. Next, suppose that  $Y$  is a constant. Then  $H(Y) = 0$ . Moreover,  $Y \preceq X_1 \wedge \dots \wedge X_n$  by definition of  $\wedge$ . Thus, by Lemma 3.3.1(c),  $H(Y|X_1 \wedge \dots \wedge X_n) = 0$ , and

$$\begin{aligned} I_\lambda(X_1, \dots, X_n : Y) &= I(X_1 \wedge \dots \wedge X_n : Y) \\ &= I(Y : X_1 \wedge \dots \wedge X_n) \\ &= H(Y) - H(Y|X_1 \wedge \dots \wedge X_n) \\ &= 0. \end{aligned}$$

$(\mathbf{Eq})$ : Consider  $X_1 \wedge \dots \wedge X_n$ . The equivalence class (with respect to  $\cong$ ) in which this random variable resides is closed under substitution of  $X_i$  (for  $i = 1, \dots, n$ ) or  $Y$  by an informationally

equivalent random variable. Hence, because

$$\begin{aligned} I_{\wedge}(X_1, \dots, X_n : Y) &= H(Y) - H(Y|X_1 \wedge \dots \wedge X_n) \\ &= H(X_1 \wedge \dots \wedge X_n) - H(X_1 \wedge \dots \wedge X_n|Y), \end{aligned}$$

by Lemma 3.3.1(a), the desired result holds.

(**TM**): Suppose that  $Y \preceq Z$ . For simplicity, let  $Q = X_1 \wedge \dots \wedge X_n$ . Then,

$$\begin{aligned} I_{\wedge}(X_1, \dots, X_n : Y) &= H(Q) - H(Q|Y) \\ &\leq H(Q) - H(Q|Z) \quad \text{by Lemma 3.3.1(b)} \\ &= I_{\wedge}(X_1, \dots, X_n : Z), \end{aligned}$$

as desired.

(**M<sub>0</sub>**): By definition of  $\wedge$ , we have  $X_1 \wedge \dots \wedge X_n \wedge W \preceq X_1 \wedge \dots \wedge X_n$ . Hence,

$$\begin{aligned} I_{\wedge}(X_1, \dots, X_n, W : Y) &= H(X_1 \wedge \dots \wedge X_n \wedge W) - H(X_1 \wedge \dots \wedge X_n \wedge W|Y) \\ &\leq H(X_1 \wedge \dots \wedge X_n) - H(X_1 \wedge \dots \wedge X_n|Y) \quad \text{by Lemma 3.3.1(b)} \\ &= I_{\wedge}(X_1, \dots, X_n : Y), \end{aligned}$$

as desired.

Next, suppose that there exists  $Z \in \{X_1, \dots, X_n\}$  such that  $Z \preceq W$ . Then, by the algebraic laws of  $\wedge$ , we have  $X_1 \wedge \dots \wedge X_n \wedge W \cong X_1 \wedge \dots \wedge X_n$ . Hence,

$$\begin{aligned} I_{\wedge}(X_1, \dots, X_n, W : Y) &= H(X_1 \wedge \dots \wedge X_n \wedge W) - H(X_1 \wedge \dots \wedge X_n \wedge W|Y) \\ &= H(X_1 \wedge \dots \wedge X_n) - H(X_1 \wedge \dots \wedge X_n|Y) \quad \text{by Lemma 3.3.1(a)} \\ &= I_{\wedge}(X_1, \dots, X_n : Y), \end{aligned}$$

as desired.

(**S<sub>0</sub>**): By the commutativity law,  $X_1 \wedge \dots \wedge X_n$  is invariant (with respect to  $\cong$ ) under reordering of  $X_1, \dots, X_n$ . Hence, the desired result follows immediately from the identity  $I_{\wedge}(X_1, \dots, X_n : Y) = H(X_1 \wedge \dots \wedge X_n) - H(X_1 \wedge \dots \wedge X_n|Y)$  and Lemma 3.3.1(a).

(**LP<sub>0</sub>**): A counterexample is provided by IMPERFECTRDN (Figure 3.5).

□

**Lemma 3.C.8.** *With respect to mutual information, the measure of intersection information  $I_{\wedge}(X_1, \dots, X_n : Y)$  satisfies (**LB**) and (**SR**), but not (**Id**).*

*Proof.* **(LB)**: Suppose that  $Q \preceq X_i$  for  $i = 1, \dots, n$ . Because  $X_1 \wedge \dots \wedge X_n$  is the largest (informationally richest) random variable that is informationally poorer than  $X_i$  for  $i = 1, \dots, n$ , it follows that  $Q \preceq X_1 \wedge \dots \wedge X_n$ . Therefore,

$$\begin{aligned} I_\wedge(X_1, \dots, X_n : Y) &= H(X_1 \wedge \dots \wedge X_n) - H(X_1 \wedge \dots \wedge X_n | Y) \\ &\geq H(Q) - H(Q | Y) \quad \text{by Lemma 3.3.1(b)} \\ &= I(Q : Y), \end{aligned}$$

as desired.

**(SR)**: By definition,  $I_\wedge(X_1 : Y) = I(X_1 : Y)$ .

**(Id)**: We have  $X \wedge Y \preceq X \vee Y$  by definition of  $\wedge$  and  $\vee$ . Thus,

$$\begin{aligned} I_\wedge(X, Y : X \vee Y) &= I(X \wedge Y : X \vee Y) \\ &= H(X \wedge Y) - H(X \wedge Y | X \vee Y) \\ &= H(X \wedge Y) \quad \text{by Lemma 3.3.1(a)} \\ &= I^0(X : Y) \\ &\neq I(X : Y) . \end{aligned}$$

□

**Lemma 3.C.9.** *The measure of intersection information  $I_\wedge(X_1, \dots, X_n : Y)$  does not satisfy **(M<sub>1</sub>)**, **(S<sub>1</sub>)**, and **(LP<sub>1</sub>)**.*

*Proof.* **(M<sub>1</sub>)**: A counterexample is provided in IMPERFECTRDN (Figure 3.5), where  $I_\wedge(X_1 : Y) = 0.99$  bits, yet  $I_\wedge(X_1, Y : Y) = 0$  bits.

**(S<sub>1</sub>)**: A counterexample. We show  $I_\wedge(X, X : Y) \neq I_\wedge(X, Y : X)$ .

$$\begin{aligned} I_\wedge(X, X : Y) - I_\wedge(X, Y : X) &= I(X : Y) - I_\wedge(X, Y : X) \\ &= I(X : Y) - I(X \wedge Y : X) \\ &= I(X : Y) - H(X \wedge Y) - H(X \wedge Y | X) \\ &= I(X : Y) - H(X \wedge Y) \\ &\neq 0 . \end{aligned}$$

**(LP<sub>1</sub>)**: This follows from not satisfying **(LP<sub>0</sub>)**.

□

### 3.D Miscellaneous Results

#### Simplification of $I_\lambda^0$

**Lemma 3.D.1.** *We have  $I_\lambda^0(X_1, \dots, X_n : Y) = H(X_1 \wedge \dots \wedge X_n \wedge Y)$ .*

*Proof.* By definition,

$$\begin{aligned} I_\lambda^0(X_1, \dots, X_n : Y) &\equiv \max_{\Pr(Q|Y)} I^0(Q : Y) \\ &\quad \text{subject to } Q \preceq X_i \ \forall i \in \{1, \dots, n\} \\ &= \max_{\Pr(Q|Y)} H(Q \wedge Y) \\ &\quad \text{subject to } Q \preceq X_i \ \forall i \in \{1, \dots, n\} \end{aligned}$$

Let  $Q$  be an arbitrary random variable satisfying the constraint  $Q \preceq X_i$  for  $i = 1, \dots, n$ . Because  $X_1 \wedge \dots \wedge X_n$  is the largest random variable (in the sense of the partial order  $\preceq$ ) that is informationally poorer than  $X_i$  for  $i = 1, \dots, n$ , we have  $Q \preceq X_1 \wedge \dots \wedge X_n$ . By the property of  $\wedge$  pointed out before, we also have  $Q \wedge Y \preceq X_1 \wedge \dots \wedge X_n \wedge Y$ . By Lemma 3.3.1(b), this implies that  $H(Q \wedge Y) \leq H(X_1 \wedge \dots \wedge X_n \wedge Y)$ . Therefore,  $I_\lambda^0(X_1, \dots, X_n : Y) = H(X_1 \wedge \dots \wedge X_n \wedge Y)$ .  $\square$

#### Simplification of $I_\lambda$

**Lemma 3.D.2.** *We have  $I_\lambda(X_1, \dots, X_n : Y) = I(X_1 \wedge \dots \wedge X_n : Y)$ .*

*Proof.* By definition,

$$\begin{aligned} I_\lambda(X_1, \dots, X_n : Y) &\equiv \max_{\Pr(Q|Y)} I(Q : Y) \\ &\quad \text{subject to } Q \preceq X_i \ \forall i \in \{1, \dots, n\} \\ &= H(Y) - \min_{\Pr(Q|Y)} H(Y|Q) \\ &\quad \text{subject to } Q \preceq X_i \ \forall i \in \{1, \dots, n\} \end{aligned}$$

Let  $Q$  be an arbitrary random variable satisfying the constraint  $Q \preceq X_i$  for  $i = 1, \dots, n$ . Because  $X_1 \wedge \dots \wedge X_n$  is the largest random variable (in the sense of the partial order  $\preceq$ ) that is informationally poorer than  $X_i$  for  $i = 1, \dots, n$ , we have  $Q \preceq X_1 \wedge \dots \wedge X_n$ . By Lemma 3.3.1(b), this implies that  $H(Y|Q) \geq H(Y|X_1 \wedge \dots \wedge X_n \wedge Y)$ . Therefore,  $I_\lambda(X_1, \dots, X_n : Y) = I(X_1 \wedge \dots \wedge X_n : Y)$ .  $\square$

**Proof that  $I_\lambda(X_1, \dots, X_n : Y) \leq I_{\min}(X_1, \dots, X_n : Y)$**

**Lemma 3.D.3.** *We have  $I_\lambda(X_1, \dots, X_n : Y) \leq I_{\min}(X_1, \dots, X_n : Y)$*



*Proof.* Starting from the definitions,

$$\begin{aligned} I_{\wedge}(X_1, \dots, X_n : Y) &\equiv I(X_1 \wedge \dots \wedge X_n : Y) \\ &= \sum_y \Pr(y) I(X_1 \wedge \dots \wedge X_n : y) \\ I_{\min}(\{X_1, \dots, X_n\} : Y) &\equiv \sum_y \Pr(y) \min_i I(X_i : y) . \end{aligned}$$

For a particular state  $y$ , without loss of generality we define the minimizing predictor  $X_m$  by  $X_m \equiv \operatorname{argmin}_{X_i} I(X_i : y)$  and the common random variable  $Q \equiv X_1 \wedge \dots \wedge X_n$ . It then remains to show that  $I(Q : y) \leq I(X_m : y)$ .

By definition of  $\wedge$ , we have  $Q \preceq X_m$ . Hence,

$$\begin{aligned} I(X_m : y) &= H(X_m) - H(X_m | Y = y) \\ &\geq H(Q) - H(Q | Y = y) \quad \text{by Lemma 3.3.1(b)} \\ &= I(Q : y) . \end{aligned}$$

□

### State-dependent zero-error information

We define the state-dependent zero-error information,  $I^0(X : Y = y)$  as,

$$I^0(X : Y = y) \equiv \log \frac{1}{\Pr(Q = q)} ,$$

where the random variable  $Q \equiv X \wedge Y$  and  $\Pr(Q = q)$  is the probability of the connected component containing state  $y \in Y$ . This entails that  $\Pr(y) \leq \Pr(q) \leq 1$ . Similar to the state-dependent information,  $\mathbb{E}_Y I^0(X : y) = I^0(X : Y)$ , where  $\mathbb{E}_Y$  is the expectation value over  $Y$ .

*Proof.* We define two functions  $f$  and  $g$ :

- $f : y \rightarrow q$  s.t.  $\Pr(q|y) = 1$  where  $q \in Q$  and  $y \in Y$ .
- $g : q \rightarrow \{y_1, \dots, y_k\}$  s.t.  $\Pr(q|y_i) = 1$  where  $q \in Q$  and  $y \in Y$ .

Now we have,

$$\mathbb{E}_Y I^0(X : y) \equiv \sum_{y \in Y} \Pr(y) \log \frac{1}{\Pr(f(y))} .$$

Since each  $y$  is associated with exactly one  $q$ , we can reindex the  $\sum_{y \in Y}$ . We then simplify to

achieve the result.

$$\begin{aligned}
\sum_{y \in Y} \Pr(y) \log \frac{1}{\Pr(f(y))} &= \sum_{q \in Q} \sum_{y \in g(q)} \Pr(y) \log \frac{1}{\Pr(f(y))} \\
&= \sum_{q \in Q} \sum_{y \in g(q)} \Pr(y) \log \frac{1}{\Pr(q)} = \sum_{q \in Q} \log \frac{1}{\Pr(q)} \sum_{y \in g(q)} \Pr(y) \\
&= \sum_{q \in Q} \log \frac{1}{\Pr(q)} \Pr(q) = \sum_{q \in Q} \Pr(q) \log \frac{1}{\Pr(q)} \\
&= H(Q) = I^0(X:Y) .
\end{aligned}$$

□

### 3.E Misc Figures

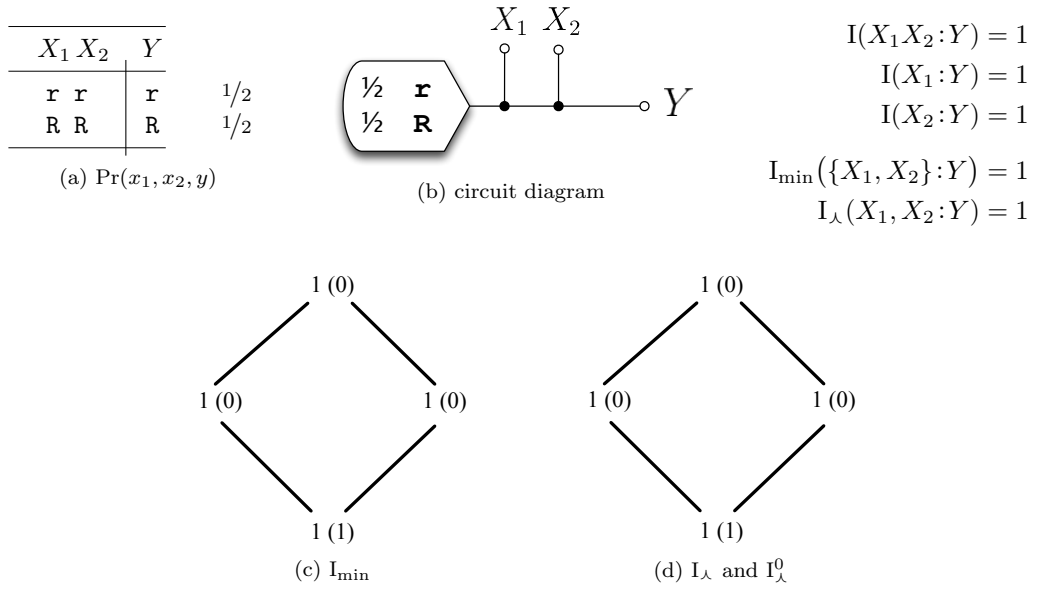


Figure 3.7: Example RDN. In this example  $I_{\min}$  and  $I_{\lambda}$  reach the same answer yet diverge drastically for example IMPERFECTRDN.

## Chapter 4

# Irreducibility is Minimum Synergy among Parts

In this chapter we explore how a collective action can be “irreducible to the actions performed by its parts”. First, we show that computing synergy among four common notions of “parts” gives rise to a spectrum of four distinct measures of irreducibility. Second, using Partial Information Decomposition[36], we introduce a nonnegative expression for each notion of irreducibility. Third, we delineate these four notions of irreducibility with exemplary binary circuits.

### 4.1 Introduction

Before we discussed computing synergy among random variables. Now we show that we can define broader notions of irreducibility by computing synergy among joint random variables. Therefore, a measure of synergy will allow us to quantify a myriad of notions of irreducibility. One pertinent application of quantifying irreducibility is finding the most useful granularity for analyzing a complex system in which interactions occur at multiple scales. Prior work [6, 19, 1, 36] has proposed measures of irreducibility, but there remains no consensus which measure is most valid.

#### 4.1.1 Notation

In our treatment of irreducibility, the  $n$  agents are random variables  $\{X_1, \dots, X_n\}$ , and the collective action the agents perform is predicting (having mutual information about) a single target random variable  $Y$ . We use the following notation throughout. Let

**X**: The set of  $n$  elementary random variables (r.v.).  $\mathbf{X} \equiv \{X_1, X_2, \dots, X_n\}$ .  $n \geq 2$ .

$X_{1\dots n}$ : The *whole*, the joint r.v. (cartesian product) of all  $n$  elements,  $X_{1\dots n} \equiv X_1 \vee \dots \vee X_n$ .

$Y$ : The “target” random variable to be predicted.

$\mathcal{P}(\mathbf{X})$ : The set of all parts (random variables) derivable from a nonempty, proper subset of  $\mathbf{X}$ . For a set of  $n$  elements there are  $2^n - 2$  possible parts. Formally,  $\mathcal{P}(\mathbf{X}) \equiv \left\{ S_1 \curlyvee \cdots \curlyvee S_{|\mathbf{S}|} : \mathbf{S} \subset \mathbf{X}, \mathbf{S} \neq \emptyset \right\}$ .

$\mathbf{P}$ : A set of  $m$  parts  $\mathbf{P} \equiv \{P_1, P_2, \dots, P_m\}$ ,  $2 \leq m \leq n$ . Each part  $P_i$  is an element (random variable) in the set  $\mathcal{P}(\mathbf{X})$ . The joint random variable of all  $m$  parts is always informationally equivalent to  $X_{1\dots n}$ , i.e.,  $P_1 \curlyvee \cdots \curlyvee P_m \cong X_{1\dots n}$ . Hereafter, the terms “part” and “component” are used interchangeably.

$A_i$ : The  $i$ ’th “Almost”. An “Almost” is a part (joint random variable) only lacking the element  $X_i$ .  $1 \leq i \leq n$ . Formally,  $A_i \equiv X_1 \curlyvee \cdots \curlyvee X_{i-1} \curlyvee X_{i+1} \curlyvee \cdots \curlyvee X_n$ .

All capital letters are random variables. All bolded capital letters are sets of random variables.

## 4.2 Four common notions of irreducibility

Prior literature [11, 1, 6, 29] has intuitively conceptualized the irreducibility of the information a whole  $X_{1\dots n}$  conveys about  $Y$  in terms of how much information about  $Y$  is lost upon “breaking up”  $X_{1\dots n}$  into a set of parts  $\mathbf{P}$ . We express this intuition formally by computing the aggregate information  $\mathbf{P}$  has about  $Y$ , and then subtracting it from the mutual information  $I(X_{1\dots n}:Y)$ . But what are the parts  $\mathbf{P}$ ? The four most common choices are:

1. **The singleton elements.** We take the set of  $n$  elements,  $\mathbf{X}$ , compute the mutual information with  $Y$  when all  $n$  elements work separately, and then subtract it from  $I(X_{1\dots n}:Y)$ . Information beyond the Elements (**lbE**) is the weakest notion of irreducibility. In the PI-diagram[36] of  $I(X_{1\dots n}:Y)$ , **lbE** is the sum of all synergistic PI-regions.
2. **Any partition of (disjoint) parts.** We enumerate all possible partitions of set  $\mathbf{X}$ . Formally, a partition  $\mathbf{P}$  is any set of parts  $\{P_1, \dots, P_m\}$  such that,  $P_i \wedge P_j \prec X_k$  where  $i, j \in \{1, \dots, m\}$ ,  $i \neq j$ , and  $k \in \{1, \dots, n\}$ . For each partition, we compute the mutual information with  $Y$  when its  $m$  parts work separately. We then take the maximum information over all partitions and subtract it from  $I(X_{1\dots n}:Y)$ . Information beyond the Disjoint Parts (**lbDp**) quantifies  $I(X_{1\dots n}:Y)$ ’s irreducibility to information conveyed by disjoint parts.
3. **Any two parts.** We enumerate all “part-pairs” of set  $\mathbf{X}$ . Formally, a part-pair  $\mathbf{P}$  is any set of exactly two elements in  $\mathcal{P}(\mathbf{X})$ . For each part-pair, we compute the mutual information with  $Y$  when the parts work separately. We then take the maximum mutual information over all part-pairs and subtract it from  $I(X_{1\dots n}:Y)$ . Information beyond the Two Parts (**lb2p**) quantifies  $I(X_{1\dots n}:Y)$ ’s irreducibility to information conveyed by any pair of parts.

4. **All possible parts.** We take the set of all possible parts of set  $\mathbf{X}$ ,  $\mathcal{P}(\mathbf{X})$ , and compute the information about  $Y$  conveyed when all  $2^n - 2$  parts work separately and subtract it from  $I(X_{1\dots n}:Y)$ . Information beyond All Parts (**lbAp**) is the strongest notion of irreducibility. In the PI-diagram of  $I(X_{1\dots n}:Y)$ , **lbAp** is the value of PI-region  $\{1\dots n\}$ .

### 4.3 Quantifying the four notions of irreducibility

To calculate the information in the whole beyond its elements, the first thing that comes to mind is to take the whole and subtract the sum over the elements, i.e.,  $I(X_{1\dots n}:Y) - \sum_{i=1}^n I(X_i:Y)$ . However, the sum *double-counts* when over multiple elements convey the same information about  $Y$ . To avoid double-counting the same information, we need to change the sum to “union”. Whereas summing adds duplicate information multiple times, unioning adds duplicate information only once. This guiding intuition of “whole minus union” leads to the definition of irreducibility as the information conveyed by the whole minus the “union information” over its parts.

We provide expressions for **lbE**, **lbDp**, **lb2p**, and **lbAp** for arbitrary  $n$ . All four equations are the information conveyed by the whole,  $I(X_{1\dots n}:Y)$ , minus the maximum union information about  $Y$  over some parts  $\mathbf{P}$ ,  $I_{\cup}(P_1, \dots, P_m:Y)$ . In PID,  $I_{\cup}$  is the dual to  $I_{\cap}$ ; they are related by the inclusion–exclusion principle. Thus if we only have a  $I_{\cap}$  measure we can always express the  $I_{\cup}$  by,

$$I_{\cup}(P_1, \dots, P_m:Y) = \sum_{\mathbf{S} \subseteq \{P_1, \dots, P_m\}} (-1)^{|\mathbf{S}|+1} I_{\cap}(S_1, \dots, S_{|\mathbf{S}|}:Y) .$$

There are currently several candidate definitions of the union information[15, 14], but for our measures to work all that is required is that the  $I_{\cup}$  measure satisfy:

- (**GP**) Global Positivity:  $I_{\cup}(P_1, \dots, P_m:Y) \geq 0$ , and  $I_{\cup}(P_1, \dots, P_m:Y) = 0$  if  $Y$  is a constant.
- (**Eq**) Equivalence-Class Invariance:  $I_{\cup}(P_1, \dots, P_m:Y)$  is invariant under substitution of  $P_i$  (for any  $i = 1, \dots, m$ ) or  $Y$  by an informationally equivalent random variable.
- (**M<sub>0</sub>**) Weak Monotonicity:  $I_{\cup}(P_1, \dots, P_m, W:Y) \geq I_{\cup}(P_1, \dots, P_m:Y)$  with equality if there exists  $P_i \in \{P_1, \dots, P_m\}$  such that  $W \preceq P_i$ .
- (**S<sub>0</sub>**) Weak Symmetry:  $I_{\cup}(P_1, \dots, P_m:Y)$  is invariant under reordering of  $P_1, \dots, P_m$ .
- (**SR**) Self-Redundancy:  $I_{\cup}(P_1:Y) = I(P_1:Y)$ . The union information a single part  $P_1$  conveys about the target  $Y$  is equal to the mutual information between  $P_1$  and the target.
- (**UB**) Upperbound:  $I_{\cup}(P_1, \dots, P_m:Y) \leq I(P_1 \vee \dots \vee P_m:Y)$ . In this particular case, the joint r.v.  $P_1 \vee \dots \vee P_m \cong X_{1\dots n}$ , so this equates to  $I_{\cup}(P_1, \dots, P_m:Y) \leq I(X_{1\dots n}:Y)$ .

### 4.3.1 Information beyond the Elements

Information beyond the Elements,  $\text{lbE}(\mathbf{X} : Y)$  quantifies how much information in  $I(X_{1\dots n} : Y)$  isn't conveyed by any element  $X_i$  for  $i \in \{1, \dots, n\}$ . The Information beyond the Elements is,

$$\text{lbE}(\mathbf{X} : Y) \equiv I(X_{1\dots n} : Y) - I_{\cup}(X_1, \dots, X_n : Y) . \quad (4.1)$$

Information beyond the Elements, or *synergistic mutual information*[15], quantifies the amount of information in  $I(X_{1\dots n} : Y)$  that *only coalitions* of elements convey.

### 4.3.2 Information beyond Disjoint Parts: $\text{lbDp}$

Information beyond Disjoint Parts,  $\text{lbDp}(\mathbf{X} : Y)$ , quantifies how much information in  $I(X_{1\dots n} : Y)$  isn't conveyed by any partition of set  $\mathbf{X}$ . Like  $\text{lbE}$ ,  $\text{lbDp}$  is the total information minus the “union information” over components. Unlike  $\text{lbE}$ , the components are not the  $n$  elements but the parts of a partition. Some algebra proves that the partition with the maximum mutual information will always be a bipartition; thus we can safely restrict the maximization to bipartitions.<sup>1</sup> Therefore to quantify  $I(X_{1\dots n} : Y)$ 's irreducibility to disjoint parts, we maximize over all  $2^{n-1} - 1$  bipartitions of set  $\mathbf{X}$ . Altogether, the Information beyond Disjoint Parts is,

$$\text{lbDp}(\mathbf{X} : Y) \equiv I(X_{1\dots n} : Y) - \max_{\substack{P_1 \in \mathcal{P}(\mathbf{X}) \\ \vdots \\ P_m \in \mathcal{P}(\mathbf{X}) \\ P_i \wedge P_j \prec X_k, \forall i \neq j \ k \in \{1, \dots, n\}}} I_{\cup}(P_1, \dots, P_m : Y) \quad (4.2)$$

$$= I(X_{1\dots n} : Y) - \max_{S \in \mathcal{P}(\mathbf{X})} I_{\cup}(S, \mathbf{X} \setminus S : Y) . \quad (4.3)$$

### 4.3.3 Information beyond Two Parts: $\text{lb2p}$

Information beyond Two Parts,  $\text{lb2p}(\mathbf{X} : Y)$ , quantifies how much information in  $I(X_{1\dots n} : Y)$  isn't conveyed by any pair of parts. Like  $\text{lbDp}$ ,  $\text{lb2p}$  subtracts the maximum union information over two parts. Unlike  $\text{lbDp}$ , the two parts aren't disjoint. Some algebra proves that the part-pair conveying the most information about  $Y$  will always be a pair of “Almosts”.<sup>2</sup> Thus, we can safely restrict the maximization over all pairs of Almosts, and we maximize over the  $\binom{n}{2} = \frac{n(n-1)}{2}$  pairs of Almosts. Altogether, the Information beyond Two Parts is,

$$\text{lb2p}(X_1, \dots, X_n : Y) \equiv I(X_{1\dots n} : Y) - \max_{\substack{P_1 \in \mathcal{P}(\mathbf{X}) \\ P_2 \in \mathcal{P}(\mathbf{X})}} I_{\cup}(P_1, P_2 : Y) \quad (4.4)$$

$$= I(X_{1\dots n} : Y) - \max_{\substack{i, j \in \{1, \dots, n\} \\ i \neq j}} I_{\cup}(A_i, A_j : Y) . \quad (4.5)$$

---

<sup>1</sup>See Appendix 4.B.1 for a proof.

<sup>2</sup>See Appendix 4.B.2 for a proof.

#### 4.3.4 Information beyond All Parts: **lbAp**

Information beyond All Parts,  $\text{lbAp}(\mathbf{X} : Y)$ , quantifies how much information in  $I(X_{1\dots n} : Y)$  isn't conveyed by any part. Like  $\text{lb2p}$ ,  $\text{lbAp}$  subtracts the union information over overlapping parts. Unlike  $\text{lb2p}$ , the union is not over two parts, but all possible parts. Some algebra proves that the entirety of the information conveyed by all  $2^n - 2$  parts working separately is equally conveyed by the  $n$  Almosts working separately.<sup>3</sup> Thus we can safely contract the union information to the  $n$  Almosts. Altogether, the Information beyond All Parts is,

$$\begin{aligned} \text{lbAp}(X_1, \dots, X_n : Y) &\equiv I(X_{1\dots n} : Y) - I_{\cup}(\mathcal{P}(\mathbf{X}) : Y) \\ &= I(X_{1\dots n} : Y) - I_{\cup}(A_1, A_2, \dots, A_n : Y) . \end{aligned} \quad (4.6)$$

Whereas Information beyond the Elements quantifies the amount of information in  $I(X_{1\dots n} : Y)$  only conveyed by coalitions, Information beyond All Parts, or *holistic mutual information*, quantifies the amount of information in  $I(X_{1\dots n} : Y)$  only conveyed by the whole.

By properties **(GP)** and **(UB)**, our four measures are nonnegative and bounded by  $I(X_{1\dots n} : Y)$ . Finally, each succeeding of notion of components is a generalization of the prior. This successive generality gives rise to the handy inequality:

$$\text{lbAp}(\mathbf{X} : Y) \leq \text{lb2p}(\mathbf{X} : Y) \leq \text{lbDp}(\mathbf{X} : Y) \leq \text{lbE}(\mathbf{X} : Y) . \quad (4.7)$$

### 4.4 Exemplary Binary Circuits

For  $n = 2$ , all four notions of irreducibility are equivalent; each one is simply the value of PI-region  $\{12\}$  (see subfigures 4.2a–d). The canonical example of irreducibility for  $n = 2$  is example XOR (Figure 4.1). In XOR, the irreducibility of  $X_1$  and  $X_2$  specifying  $Y$  is analogous to irreducibility of hydrogen and oxygen extinguishing fire. The whole  $X_1X_2$  fully specifies  $Y$ ,  $I(X_1X_2 : Y) = H(Y) = 1$  bit, but  $X_1$  and  $X_2$  separately convey nothing about  $Y$ ,  $I(X_1 : Y) = I(X_2 : Y) = 0$  bits.

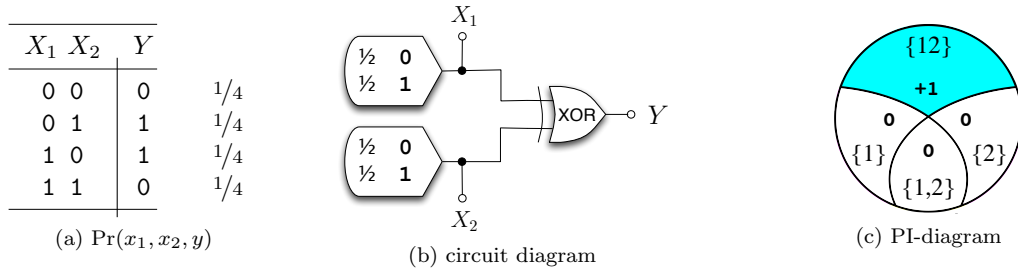


Figure 4.1: Example XOR.  $X_1X_2$  irreducibly specifies  $Y$ .  $I(X_1X_2 : Y) = H(Y) = 1$  bit.

<sup>3</sup>See Appendix 4.B.3 for a proof.

For  $n > 2$ , the four notions of irreducibility diverge; subfigures 4.2e–j depicts **lbE**, **lbAp**, **lbDp**, and **lb2p** when  $n = 3$ . We provide exemplary binary circuits delineating each measure. Every circuit has  $n = 3$  elements, meaning  $\mathbf{X} = \{X_1, X_2, X_3\}$ , and build atop example XOR.

#### 4.4.1 XorUnique: Irreducible to elements, yet reducible to a partition

To concretize how a collective action could be irreducible to elements, yet still reducible to a partition, consider a hypothetical set of agents  $\{X_1, X_2, \dots, X_{100}\}$  where the first 99 agents cooperate to specify  $Y$ , but agent  $X_{100}$  doesn't cooperate with the joint random variable  $X_1 \cdots X_{99}$ . The **lbE** among these 100 agents would be *positive*, however, **lbDp** would be *zero* because the work that  $X_1 \cdots X_{100}$  performs can be reduced to two disjoint parts,  $X_1 \cdots X_{99}$  and  $X_{100}$ , working separately.

Example XORUNIQUE (Figure 4.3) is analogous to the situation above. The whole specifies two bits of uncertainty,  $I(X_1 X_2 X_3 : Y) = H(Y) = 2$  bits. The doublet  $X_1 X_2$  solely specifies the “digit-bit” of  $Y$  (0/1),  $I(X_1 X_2 : Y) = 1$  bit, and the singleton  $X_3$  solely specifies the “letter-bit” of  $Y$  (a/A),  $I(X_3 : Y) = 1$  bit. We apply each notion of irreducibility to XORUNIQUE:

**lbE** How much of  $X_1 X_2 X_3$ 's information about  $Y$  can be reduced to the information conveyed by the singleton elements working separately? Working alone,  $X_3$  still specifies the letter-bit of  $Y$ , but  $X_1$  nor  $X_2$  can unilaterally specify the digit-bit of  $Y$ ,  $I(X_1 : Y) = 0$  and  $I(X_2 : Y) = 0$  bits. As only the letter-bit is specified when the three singletons work separately,  $\text{lbE}(\mathbf{X} : Y) = I(X_1 X_2 X_3 : Y) - 1 = 2 - 1 = 1$  bit.

**lbDp** How much of  $X_1 X_2 X_3$ 's information about  $Y$  can be reduced to the information conveyed by disjoint parts working separately? Per subfigures 4.2g–i, there are three bipartitions of  $X_1 X_2 X_3$ , and one of them is  $\{X_1 X_2, X_3\}$ . The doublet part  $X_1 X_2$  specifies the digit-bit of  $Y$ , and the singleton part  $X_3$  specifies the letter-bit of  $Y$ . As there is a partition of  $X_1 X_2 X_3$  that fully accounts for  $X_1 X_2 X_3$ 's specification of  $Y$ ,  $\text{lbDp}(\mathbf{X} : Y) = 2 - 2 = 0$  bits.

**lb2p/lbAp** How much of  $X_1 X_2 X_3$ 's information about  $Y$  can be reduced to the information conveyed by two parts working separately? From above we see that **lbDp** is zero bits. Per eq. (4.7), **lb2p** and **lbAp** are stricter notions of irreducibility than **lbDp**, therefore **lb2p** and **lbAp** must also be zero bits.

#### 4.4.2 DoubleXor: Irreducible to a partition, yet reducible to a pair

In example DOUBLEXOR (Figure 4.4), the whole specifies two bits,  $I(X_1 X_2 X_3 : Y) = H(Y) = 2$  bits. The doublet  $X_1 X_2$  solely specifies the “left-bit”, and the doublet  $X_2 X_3$  solely specifies the “right-bit”. Applying each notion of irreducibility to DOUBLEXOR:



**lbE** How much of  $X_1X_2X_3$ 's information about  $Y$  can be reduced to the information conveyed by singleton elements? The three singleton elements specify nothing about  $Y$ ,  $I(X_i:Y) = 0$  bits  $\forall i$ . This means the whole is utterly irreducible to its elements, making  $\text{lbE}(\mathbf{X} : Y) = I(X_1X_2X_3:Y) - 0 = 2$  bits.

**lbDp** How much of  $X_1X_2X_3$ 's information about  $Y$  can be reduced to the information conveyed by disjoint parts? Per subfigures 4.2g-i, the three bipartitions of  $X_1X_2X_3$  are:  $\{X_1X_2, X_3\}$ ,  $\{X_1X_3, X_2\}$ , and  $\{X_2X_3, X_1\}$ . In the first bipartition,  $\{X_1X_2, X_3\}$ , the doublet  $X_1X_2$  specifies the left-bit of  $Y$  and the singleton  $X_3$  specifies nothing for a total of one bit. Similarly, in the second bipartition,  $\{X_2X_3, X_1\}$ ,  $X_2X_3$  specifies the right-bit of  $Y$  and the singleton  $X_1$  specifies nothing for a total of one bit. Finally, in the bipartition  $\{X_1X_3, X_2\}$  both  $X_1X_3$  and  $X_2$  specify nothing for a total of zero bits. Taking the maximum over the three bipartitions,  $\max[1, 1, 0] = 1$ , we discover disjoint parts specify at most one bit, leaving  $\text{lbDp}(\mathbf{X} : Y) = I(X_1X_2X_3:Y) - 1 = 2 - 1 = 1$  bit.

**lb2p** How much of  $X_1X_2X_3$ 's information about  $Y$  can be reduced to the information conveyed by two parts? Per subfigures 4.2k-j, there are three pairs of Almosts, and one of them is  $\{X_1X_2, X_1X_3\}$ . The Almost  $X_1X_2$  specifies the left-bit of  $Y$ , and the Almost  $X_1X_3$  specifies the right-bit of  $Y$ . As there is a pair of parts that fully accounts for  $X_1X_2X_3$ 's specification of  $Y$ ,  $\text{lb2p}(\mathbf{X} : Y) = 0$  bits.

**lbAp** How much of  $X_1X_2X_3$ 's information about  $Y$  can be reduced to the information conveyed by all possible parts? From above we see that  $\text{lb2p}$  is zero bits. Per eq. (4.7),  $\text{lbAp}$  is stricter than  $\text{lb2p}$ , therefore  $\text{lbAp}$  is also zero bits.

#### 4.4.3 TripleXor: Irreducible to a pair of components, yet still reducible

Example TRIPLEXOR (Figure 4.5) has trifold symmetry and the whole specifies three bits,  $I(X_1X_2X_3:Y) = H(Y) = 3$  bits. Each bit is solely specified by one of three doublets:  $X_1X_2$ ,  $X_1X_3$ , or  $X_2X_3$ . Applying each notion of irreducibility to TRIPLEXOR:

**lbE** Working individually, the three elements specify absolutely nothing about  $Y$ ,  $I(X_1:Y) = I(X_2:Y) = I(X_3:Y) = 0$  bits. Thus, the whole is utterly irreducible to elements, making  $\text{lbE}(\mathbf{X} : Y) = I(X_1X_2X_3:Y) - 0 = 3$  bits.

**lbDp** The three bipartitions of  $X_1X_2X_3$  are:  $\{X_1X_2, X_3\}$ ,  $\{X_1X_3, X_2\}$ , and  $\{X_2X_3, X_1\}$ . In the first bipartition,  $\{X_1X_2, X_3\}$ , the doublet  $X_1X_2$  specifies one bit of  $Y$  and the singleton  $X_3$  specifies nothing for a total of one bit. By TRIPLEXOR's trifold symmetry, we get the same value for bipartitions  $\{X_1X_2, X_3\}$  and  $\{X_2X_3, X_1\}$ . Taking the maximum over the three bipartitions,

$\max[1, 1, 1] = 1$ , we discover a partition specifies at most one bit, leaving  $\text{lbDp}(\mathbf{X} : Y) = I(X_1X_2X_3 : Y) - 1 = 2$  bits.

**lb2p** There are three pairs of Almosts:  $\{X_1X_2, X_2X_3\}$ ,  $\{X_1X_2, X_1X_3\}$ , and  $\{X_1X_3, X_2X_3\}$ . Each pair of Almosts specifies exactly two bits. Taking the maximum over the pairs,  $\max[2, 2, 2] = 2$ , we discover a pair of parts specifies at most two bits, leaving  $\text{lb2p}(\mathbf{X} : Y) = I(X_1X_2X_3 : Y) - 2 = 3 - 2 = 1$  bit.

**lbAp** The  $n$  Almosts of  $X_1X_2X_3$  are  $\{X_1, X_2, X_1X_3, X_2X_3\}$ . Each Almost specifies one bit of  $Y$ , for a total of three bits, making  $\text{lbAp}(\mathbf{X} : Y) = I(X_1X_2X_3 : Y) - 3 = 0$  bits.

#### 4.4.4 Parity: Complete irreducibility

In example PARITY (Figure 4.6), the whole specifies one bit of uncertainty,  $I(X_1X_2X_3 : Y) = H(Y) = 1$  bit. No singleton or doublet specifies anything about  $Y$ ,  $I(X_i : Y) = I(X_iX_j : Y) = 0$  bits  $\forall i, j$ . Applying each notion of irreducibility to PARITY:

**lbE** The whole specifies one bit, yet the elements  $\{X_1, X_2, X_3\}$  specify nothing about  $Y$ . Thus the whole is utterly irreducible to elements, making  $\text{lbE}(\mathbf{X} : Y) = I(X_1X_2X_3 : Y) - 0 = 1$  bit.

**lbDp** The three bipartitions of  $\mathbf{X}$  are:  $\{X_1X_2, X_3\}$ ,  $\{X_1X_3, X_2\}$ , and  $\{X_2X_3, X_1\}$ . By the above each doublet and singleton specifies nothing about  $Y$ , and thus each partition specifies nothing about  $Y$ . Taking the maximum over the bipartitions yields  $\max[0, 0, 0] = 0$ , making  $\text{lbDp}(\mathbf{X} : Y) = 1 - 0 = 1$  bit.

**lb2p** The pairs of  $\mathbf{X}$ 's Almosts are:  $\{X_1X_2, X_1X_3\}$ ,  $\{X_1X_2, X_2X_3\}$ , and  $\{X_1X_3, X_2X_3\}$ . As before, each doublet specifies nothing about  $Y$ , and a pair of nothings is still nothing. Taking the maximum yields  $\max[0, 0, 0] = 0$ , making  $\text{lb2p}(\mathbf{X} : Y) = 1 - 0 = 1$  bit.

**lbAp** The three Almosts of  $\mathbf{X}$  are:  $\{X_1X_2, X_1X_3, X_2X_3\}$ . Each Almost specifies nothing, and a triplet of nothings is still nothing, making  $\text{lbAp}(\mathbf{X} : Y) = 1 - 0 = 1$  bit.

Table 4.1 summarizes the results of our four irreducibility measures applied to our examples.

## 4.5 Conclusion

Within the Partial Information Decomposition framework[36], synergy is the simplest case of the broader notion of irreducibility. PI-diagrams, a generalization of Venn diagrams, are immensely helpful in improving one's intuition for synergy and irreducibility.

Example	$I(X_{1\dots n}:Y)$	lbE	lbDp	lb2p	lbAp
RDN (Fig. 1.3)	1	0	0	0	0
UNQ (Fig. 1.4)	2	0	0	0	0
XOR (Fig. 1.5)	1	1	1	1	1
XORUNIQUE (Fig. 4.3)	2	1	0	0	0
DOUBLEXOR (Fig. 4.4)	2	2	1	0	0
TRIPLEXOR (Fig. 4.5)	3	3	2	1	0
PARITY (Fig. 4.6)	1	1	1	1	1

Table 4.1: Irreducibility values for our exemplary binary circuits.

We define the irreducibility of the mutual information a set of  $n$  random variables  $\mathbf{X} = \{X_1, \dots, X_n\}$  convey about a target  $Y$  as the information the whole conveys about  $Y$ ,  $I(X_{1\dots n}:Y)$ , minus the maximum union-information conveyed by the “parts” of  $\mathbf{X}$ . The four common notions of  $\mathbf{X}$ ’s parts are: (1) the set of the  $n$  atomic elements; (2) all partitions of disjoint parts; (3) all pairs of parts; and (4) the set of all  $2^n - 2$  possible parts. All four definitions of parts are equivalent when the whole consists of two atomic elements ( $n = 2$ ), but they diverge for  $n > 2$ .

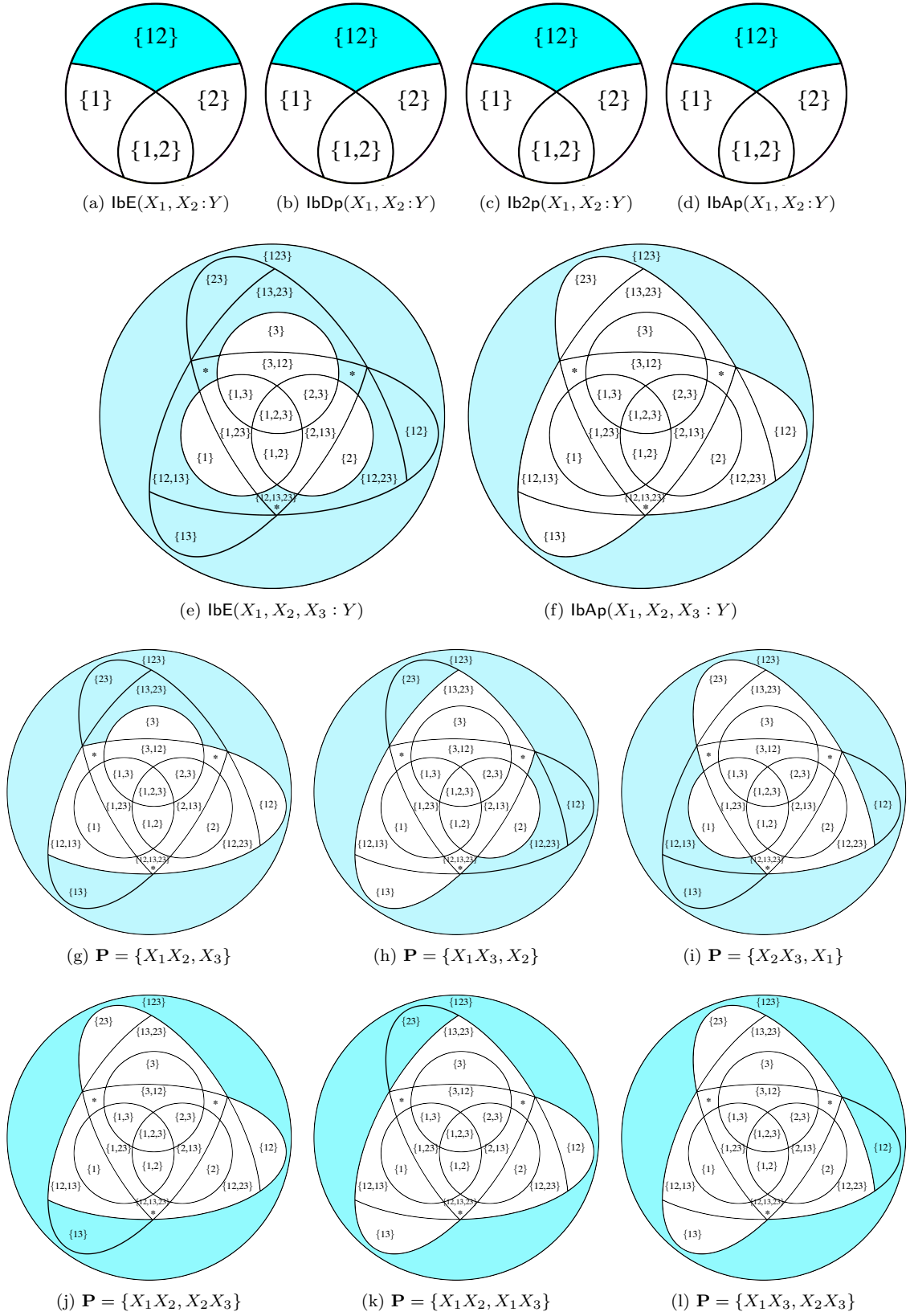


Figure 4.2: PI-diagrams depicting our four irreducibility measures when  $n = 2$  and  $n = 3$  in subfigures (a)–(d) and (e)–(l) respectively. For  $n = 3$ :  $\text{lbE}$  is (e),  $\text{lbAp}$  is (f),  $\text{lbDp}$  is the *minimum* value over subfigures (g)–(i), and  $\text{lb2p}$  is the *minimum* value over subfigures (j)–(l).

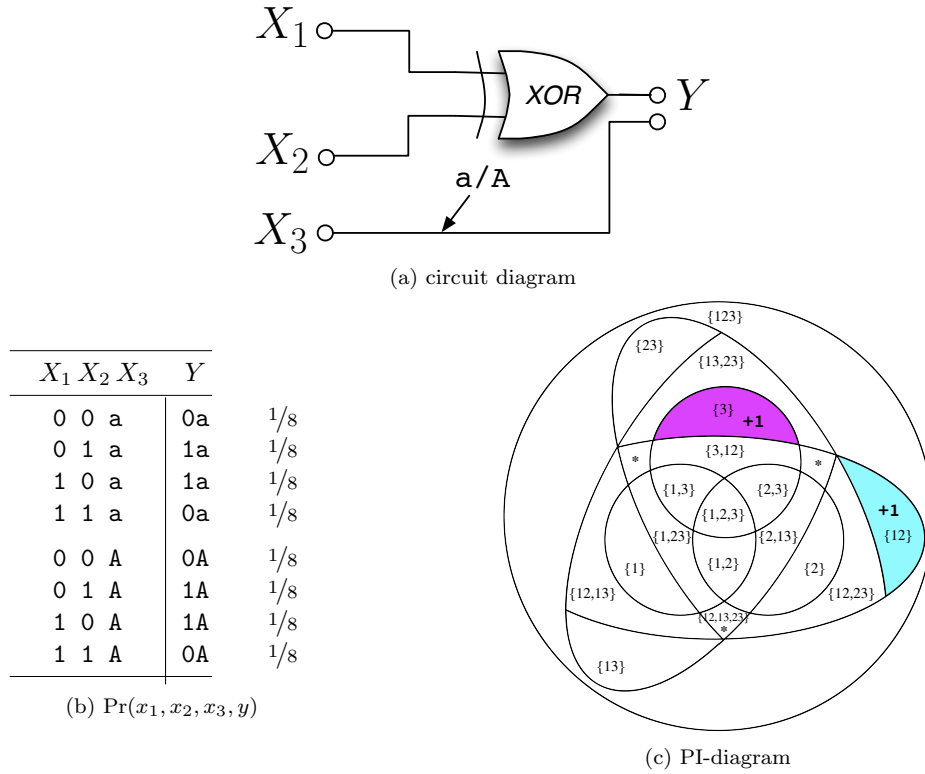


Figure 4.3: Example XORUNIQUE. Target  $Y$  has two bits of uncertainty. The doublet  $X_1X_2$  specifies the “digit bit”, and the singleton  $X_3$  specifies the “letter bit” for a total of  $I(X_1X_2:Y) + I(X_3:Y) = H(Y) = 2$  bits.  $X_1X_2X_3$ ’s specification of  $Y$  is irreducible to singletons yet fully reduces to the disjoint parts  $\{X_1X_2, X_3\}$ .

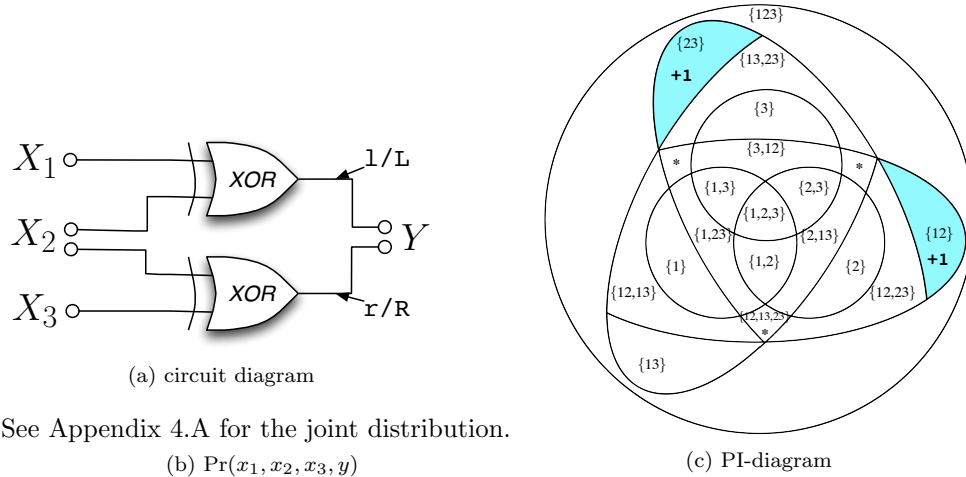
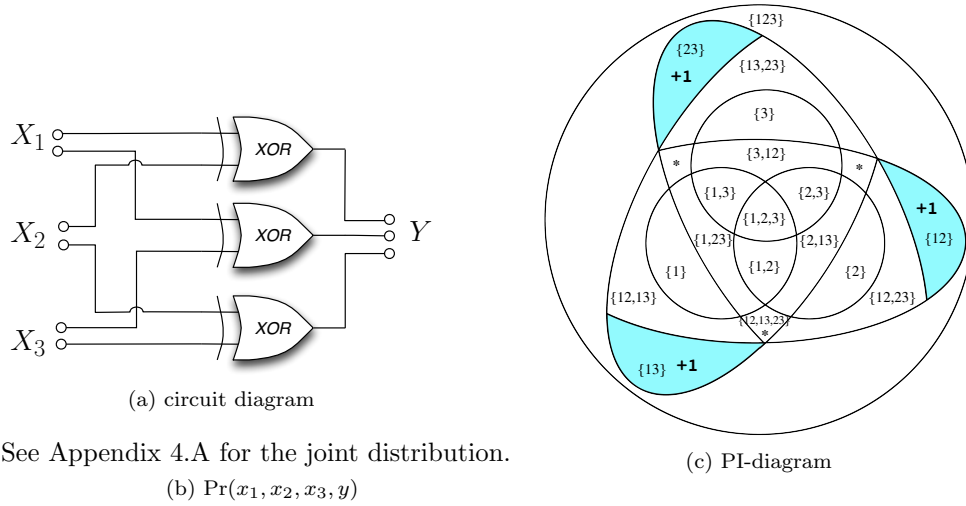


Figure 4.4: Example DOUBLEXOR. Target  $Y$  has two bits of uncertainty. The doublet  $X_1X_2$  specifies the “left bit” (1/L) and doublet  $X_2X_3$  specifies the “right bit” (r/R) for a total of  $I(X_1X_2:Y) + I(X_2X_3:Y) = H(Y) = 2$  bits.  $X_1X_2X_3$ ’s specification of  $Y$  is irreducible to disjoint parts yet fully reduces to the pair of parts  $\{X_1X_2, X_2X_3\}$ .



See Appendix 4.A for the joint distribution.

Figure 4.5: Example TRIPLEXOR. Target  $Y$  has three bits of uncertainty. Each doublet part of  $X_1X_2X_3$  specifies a distinct bit of  $Y$ , for a total of  $I(X_1X_2:Y) + I(X_1X_3:Y) + I(X_2X_3:Y) = H(Y) = 3$  bits. The whole's specification of  $Y$  is irreducible to any pair of Almosts yet fully reduces to all Almosts.

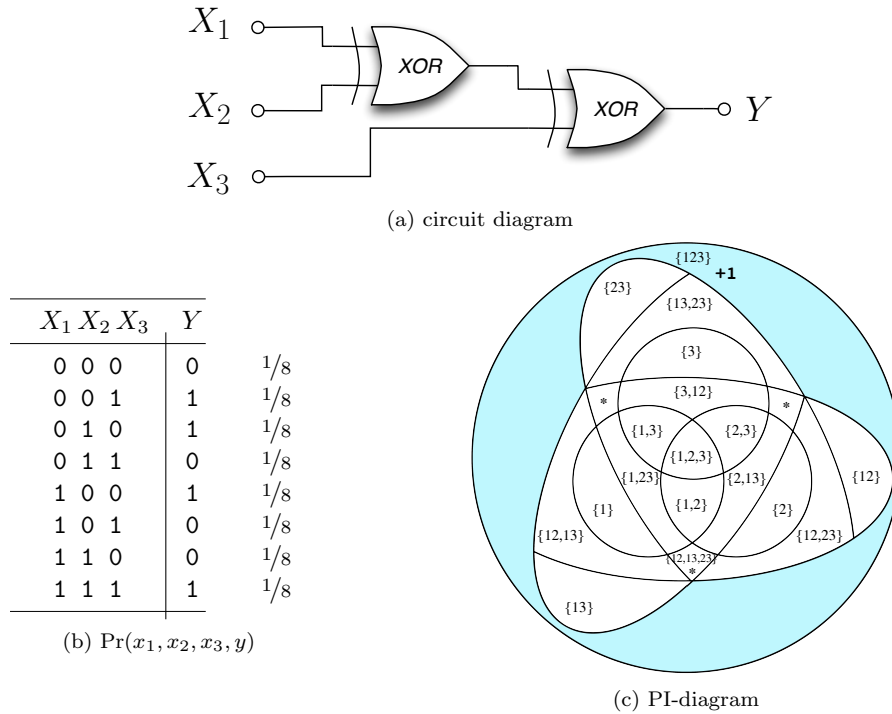


Figure 4.6: Example PARITY. Target  $Y$  has one bit of uncertainty, and only the whole specifies  $Y$ ,  $I(X_1X_2X_3:Y) = H(Y) = 1$  bit.  $X_1X_2X_3$ 's specification of  $Y$  is utterly irreducible to any collection of  $X_1X_2X_3$ 's parts, and  $\text{lbAp}(\{X_1, X_2, X_3\} : Y) = 1$  bit.

# Appendix

## 4.A Joint distributions for DoubleXor and TripleXor

$X_1$	$X_2$	$X_3$	$Y$	
0	00	0	1r	$1/16$
0	01	0	1R	$1/16$
0	10	0	Lr	$1/16$
0	11	0	LR	$1/16$
0	00	1	1R	$1/16$
0	01	1	1r	$1/16$
0	10	1	LR	$1/16$
0	11	1	Lr	$1/16$
1	00	0	Lr	$1/16$
1	01	0	LR	$1/16$
1	10	0	1r	$1/16$
1	11	0	1R	$1/16$
1	00	1	LR	$1/16$
1	01	1	Lr	$1/16$
1	10	1	1R	$1/16$
1	11	1	1r	$1/16$

Figure 4.7: Joint distribution  $\Pr(x_1, x_2, x_3, y)$  for example DOUBLEXOR.

## 4.B Proofs

**Lemma 4.B.1.** *We prove that Information beyond the Bipartition,  $lb2p(\mathbf{X} : Y)$ , equals Information beyond the Disjoint Parts,  $lbDp(\mathbf{X} : Y)$  by showing,*

$$lbDp(\mathbf{X} : Y) \leq lb2p(\mathbf{X} : Y) \leq lbDp(\mathbf{X} : Y) .$$

$X_1$	$X_2$	$X_3$	$Y$		$X_1$	$X_2$	$X_3$	$Y$	
00	00	00	000	$1/64$	10	00	00	110	$1/64$
00	00	01	001	$1/64$	10	00	01	111	$1/64$
00	00	10	010	$1/64$	10	00	10	100	$1/64$
00	00	11	011	$1/64$	10	00	11	101	$1/64$
00	01	00	001	$1/64$	10	01	00	111	$1/64$
00	01	01	000	$1/64$	10	01	01	110	$1/64$
00	01	10	011	$1/64$	10	01	10	101	$1/64$
00	01	11	010	$1/64$	10	01	11	100	$1/64$
00	10	00	100	$1/64$	10	10	00	010	$1/64$
00	10	01	101	$1/64$	10	10	01	011	$1/64$
00	10	10	110	$1/64$	10	10	10	000	$1/64$
00	10	11	111	$1/64$	10	10	11	001	$1/64$
00	11	00	101	$1/64$	10	11	00	011	$1/64$
00	11	01	100	$1/64$	10	11	01	010	$1/64$
00	11	10	111	$1/64$	10	11	10	001	$1/64$
00	11	11	110	$1/64$	10	11	11	000	$1/64$
01	00	00	000	$1/64$	11	00	00	110	$1/64$
01	00	01	001	$1/64$	11	00	01	111	$1/64$
01	00	10	010	$1/64$	11	00	10	100	$1/64$
01	00	11	011	$1/64$	11	00	11	101	$1/64$
01	01	00	001	$1/64$	11	01	00	011	$1/64$
01	01	01	000	$1/64$	11	01	01	010	$1/64$
01	01	10	011	$1/64$	11	01	10	001	$1/64$
01	01	11	010	$1/64$	11	01	11	000	$1/64$
01	10	00	100	$1/64$	11	10	00	010	$1/64$
01	10	01	101	$1/64$	11	10	01	011	$1/64$
01	10	10	110	$1/64$	11	10	10	000	$1/64$
01	10	11	111	$1/64$	11	10	11	001	$1/64$
01	11	00	101	$1/64$	11	11	00	011	$1/64$
01	11	01	100	$1/64$	11	11	01	010	$1/64$
01	11	10	111	$1/64$	11	11	10	001	$1/64$
01	11	11	110	$1/64$	11	11	11	000	$1/64$

Figure 4.8: Joint distribution  $\Pr(x_1, x_2, x_3, y)$  for example TRIPLEXOR.

*Proof.* We first show that  $\text{lbDp}(\mathbf{X} : Y) \leq \text{lb2p}(\mathbf{X} : Y)$ . By their definitions:

$$\text{lbDp}(\mathbf{X} : Y) \equiv I(Y : X_{1\dots n}) - \max_{\mathbf{P}} I_{\cup}(Y : \mathbf{P}) \quad (4.8)$$

$$\text{lbB}(\mathbf{X} : Y) \equiv I(Y : X_{1\dots n}) - \max_{S \subset \mathbf{X}} I_{\cup}(Y : \{S, \mathbf{X} \setminus S\}) \quad (4.9)$$

$$= I(Y : X_{1\dots n}) - \max_{|\mathbf{P}|=2} I_{\cup}(Y : \mathbf{P}) , \quad (4.10)$$

where  $\mathbf{P}$  enumerates over all disjoint parts of  $\mathbf{X}$ .

By removing the restriction that  $|\mathbf{P}| = 2$  from the minimized union-information in  $\text{lbB}$ , we arrive



at  $\text{lbDp}$ . As removing a restriction can only decrease the minimum, therefore  
 $\text{lbDp}(\mathbf{X} : Y) \leq \text{lbB}(\mathbf{X} : Y)$ .  $\square$

We next show that  $\text{lbB}(\mathbf{X} : Y) \leq \text{lbDp}(\mathbf{X} : Y)$ . Meaning we must show that,

$$I(X_{1\dots n} : Y) - \max_{\substack{\mathbf{P} \\ |\mathbf{P}|=2}} I_{\cup}(\mathbf{P} : Y) \leq I(X_{1\dots n} : Y) - \max_{\mathbf{P}} I_{\cup}(\mathbf{P} : Y) . \quad (4.11)$$

*Proof.* By subtracting  $I(X_{1\dots n} : Y)$  from each side and multiplying each side by  $-1$  we have,

$$\max_{\substack{\mathbf{P} \\ |\mathbf{P}|=2}} I_{\cup}(\mathbf{P} : Y) \geq \max_{\mathbf{P}} I_{\cup}(\mathbf{P} : Y) . \quad (4.12)$$

Without loss of generality, we take any individual subset/part  $S$  in  $X$ . Then we have a bipartition  $\mathbf{B}$  of parts  $\{S, \mathbf{X} \setminus S\}$ . We then further partition the part  $\mathbf{X} \setminus S$  into  $k$  disjoint subcomponents denoted  $\{T_1, \dots, T_k\}$  where  $2 \leq k \leq n - |S|$ , creating an arbitrary partition  $\mathbf{P} = \{S, T_1, \dots, T_k\}$ . We now need to show that,

$$I_{\cup}(\{S, \mathbf{X} \setminus S\} : Y) \geq I_{\cup}(\{S, T_1, \dots, T_k\} : Y) . \quad (4.13)$$

By the monotonicity axiom **(M)**, we can append each subcomponent  $T_1, \dots, T_k$  to  $\mathbf{B}$  without changing the union-information, because every subcomponent  $T_i$  is a subset of the element  $\mathbf{X} \setminus S$ . Then, using the symmetry axiom **(S<sub>0</sub>)**, we re-order the parts so that  $S, T_1, \dots, T_k$  come first. This yields,

$$I_{\cup}(\{S, T_1, \dots, T_k, \mathbf{X} \setminus S\} : Y) \geq I_{\cup}(\{S, T_1, \dots, T_k\} : Y) . \quad (4.14)$$

Applying the monotonicity axiom **(M)** again, we know that adding the entry  $\mathbf{X} \setminus S$  can only increase the union information. Therefore, we prove eq. (5.13), which proves eq. (4.11).  $\square$

Finally, by the squeeze theorem we complete the proof of eq. (5.11), that  
 $\text{lbB}(\mathbf{X} : Y) = \text{lbDp}(\mathbf{X} : Y)$ .

**Lemma 4.B.2.** *Proof that pairs of Almosts cover lb2p. We prove that the maximum union-information over all possible pairs of parts  $\{P_1, P_2\}$ , equates to the maximum union-information over all pairs of Almosts  $\{A_i, A_j\}$   $i \neq j$ . Mathematically,*

$$\max_{\substack{P_1, P_2 \\ P_1, P_2 \subset \mathbf{X}}} I_{\cup}(\{P_1, P_2\} : Y) = \max_{\substack{i, j \in \{1, \dots, n\} \\ i \neq j}} I_{\cup}(\{A_i, A_j\} : Y) . \quad (4.15)$$

*Proof.* By the right-monotonicity lemma **(RM)**, the union-information can only increase when increasing the size of the parts  $P_1$  and  $P_2$ . We can therefore ignore all parts  $P_1, P_2$  of size less than  $n - 1$ ,

$$\max_{\substack{P_1, P_2 \\ P_1, P_2 \subset \mathbf{X}}} I_{\cup}(\{P_1, P_2\} : Y) = \max_{\substack{P_1, P_2 \\ P_1, P_2 \in \mathcal{P}(\mathbf{X}) \\ |P_1|=|P_2|=n-1}} I_{\cup}(\{P_1, P_2\} : Y) \quad (4.16)$$

$$= \max_{i, j \in \{1, \dots, n\}} I_{\cup}(\{A_i, A_j\} : Y) . \quad (4.17)$$

Then by the idempotency axiom **(I)** and the monotonicity axiom **(M)**, having  $i \neq j$  can only increase the union information. Therefore,

$$\max_{i, j \in \{1, \dots, n\}} I_{\cup}(\{A_i, A_j\} : Y) = \max_{\substack{i, j \in \{1, \dots, n\} \\ i \neq j}} I_{\cup}(\{A_i, A_j\} : Y) . \quad (4.18)$$

With eq. (4.18) in hand, we easily show that the Information beyond all pairs of Subsets, **lb2p**, equates to the information beyond all pairs of Almosts:

$$\text{lb2p}(\mathbf{X} : Y) \equiv I(X_{1\dots n} : Y) - \max_{\substack{P_1, P_2 \\ P_1, P_2 \in \mathcal{P}(\mathbf{X})}} I_{\cup}(\{P_1, P_2\} : Y) \quad (4.19)$$

$$= I(X_{1\dots n} : Y) - \max_{\substack{i, j \in \{1, \dots, n\} \\ i \neq j}} I_{\cup}(\{A_i, A_j\} : Y) . \quad (4.20)$$

□

**Lemma 4.B.3.** *Proof that Almosts cover lbAp.* We wish to show that the union-information over all distinct parts of  $n$  elements,  $\mathcal{P}(\mathbf{X})$ , is equivalent to the union information over the  $n$  Almosts. Mathematically,

$$I_{\cup}(\mathcal{P}(\mathbf{X}) : Y) = I_{\cup}(\{A_1, \dots, A_n\} : Y) . \quad (4.21)$$

*Proof.* Every element in the set of parts  $\mathcal{P}(\mathbf{X})$  that isn't an Almost is a subset of an Almost. Therefore, by the monotonicity axiom **(M)** we can remove this entry. Repeating this process, we remove all entries except the  $n$  Almosts. Therefore,  $I_{\cup}(\mathcal{P}(\mathbf{X}) : Y) = I_{\cup}(\{A_1, \dots, A_n\} : Y)$ . □

## Part III

# Applications

## Chapter 5

# Improving the $\phi$ Measure

### 5.1 Introduction

The measure of integrated information,  $\phi$ , is an attempt to quantify the magnitude of conscious experience. It has a long history [31, 3, 33], and at least three different measures have been called  $\phi$ . Here we consider some adjustments to the  $\phi$  measure from [3] to correct perceived deficiencies.<sup>1</sup>

The  $\phi$  measure aims to quantify a system’s “functional irreducibility to disjoint parts.” As discussed in Chapter 4, we can use Partial Information Decomposition (PID) to derive a principled measure of irreducibility to disjoint parts. This revised measure,  $\psi$ , has numerous desirable properties over  $\phi$ .

### 5.2 Preliminaries

#### 5.2.1 Notation

We use the following notation throughout this chapter:

$n$ : the number of indivisible elements in network  $X$ .  $n \geq 2$ .

$\mathbf{P}$ : a partition of the  $n$  indivisible nodes clustered into  $m$  parts. Each part has at least one node and each partition has at least two parts, so  $2 \leq m \leq n$ .

$X_i^{\mathbf{P}}$ : a random variable representing a part  $i$  at time=0.  $1 \leq i \leq m$ .

$Y_i^{\mathbf{P}}$ : a random variable representing part  $i$  after  $t$  updates.  $1 \leq i \leq m$ .

$X$ : a random variable representing the entire network at time=0.  $X \equiv X_1^{\mathbf{P}} \cdots X_m^{\mathbf{P}}$ .

---

<sup>1</sup>We chose the 2008 version [3] because it is the most recent purely information-theoretic  $\phi$ . The most recent version from [33] uses the Hamming Distance among states and thus changes depending on the chosen labels. We are aware of no other info-theoretic measure that changes under relabeling. Secondly, the measure in [33] is in units bits-squared, which has no known information-theoretic interpretation.

$Y$ : a random variable representing the entire network after  $t$  applications of the neural network's update rule.  $Y \equiv Y_1^{\mathbf{P}} \cdots Y_m^{\mathbf{P}}$ .

$y$ : a single state of the random variable  $Y$ .

$\mathbf{X}$ : The set of  $n$  indivisible elements at time=0.

For readers accustomed to the notation in [3], the translation is  $X \equiv X_0$ ,  $Y \equiv X_1$ ,  $X_i^{\mathbf{P}} \equiv M_0^i$ , and  $Y_i^{\mathbf{P}} \equiv M_1^i$ .

For pedagogical purposes we confine this paper to deterministic networks. Therefore all remaining entropy at time  $t$  conveys information about the past, i.e.,  $I(X:Y) = H(Y)$  and  $I(X:Y_i^{\mathbf{P}}) = H(Y_i^{\mathbf{P}})$  where  $I(\bullet:\bullet)$  is the mutual information and  $H(\bullet)$  is the Shannon entropy[9]. Our model generalizes to probabilistic units with any finite number of discrete—but not continuous—states[5]. All logarithms are  $\log_2$ . All calculations are in *bits*.

### 5.2.2 Model assumptions

- (A) The  $\phi$  measure is a *state-dependent* measure, meaning that every output state  $y \in Y$  has its own  $\phi$  value. To simplify cross-system comparisons, some researchers[5] prefer to consider only the averaged  $\phi$ , denoted  $\langle \phi \rangle$ . Here we adhere to the original theoretical state-dependent formulation. But when comparing large numbers of networks we use  $\langle \phi \rangle$  for convenience.
- (B) The  $\phi$  measure aims to quantify “information intrinsic to the system”. This is often thought to be synonymous with causation, but it’s not entirely clear. But for this reason, in [3] all random variables at time=0, i.e.  $X$  and  $X_1^{\mathbf{P}}, \dots, X_m^{\mathbf{P}}$  are made to follow an *independent discrete uniform distribution*. There are actually several plausible choices for the distribution on  $X$ , but for easier comparison to [3], here we also take  $X$  to be an independent discrete uniform distribution. This means that  $H(X) = \log_2 |X|$ ,  $H(X_i^{\mathbf{P}}) = \log_2 |X_i^{\mathbf{P}}|$  where  $|\bullet|$  is the number of states in the random variable, and  $I(X_i^{\mathbf{P}}:X_j^{\mathbf{P}}) = 0 \ \forall i \neq j$ .
- (C) We set  $t = 1$ , meaning we compute these informational measures for a system undergoing a single update from time=0 to time=1. This has no impact on generality (see Appendix 5.E). To analyze real biological networks one would sweep  $t$  over all reasonable timescales, choosing the  $t$  that maximizes the complexity metric.

## 5.3 How $\phi$ works

The  $\phi$  measure has four steps and proceeds as follows:

1. For a given state  $y \in Y$ , [3] first defines the state’s *effective information*, quantifying the total magnitude of information the state  $y$  conveys about  $X$ , the r.v. representing a maximally

ignorant past. This turns out to be identical to the “specific-surprise” measure,  $I(X:y)$ , from [10],

$$\mathbf{ei}(X \rightarrow y) = I(X:y) = D_{\text{KL}} \left[ \Pr(X|y) \parallel \Pr(X) \right] . \quad (5.1)$$

Given  $X$  follows a discrete uniform distribution (assumption **(B)**),  $\mathbf{ei}(X \rightarrow y)$  simplifies to,

$$\begin{aligned} \mathbf{ei}(X \rightarrow y) &= H(X) - H(X|y) \\ &= H(X) - \sum_{x \in X} \Pr(x|y) \log \frac{1}{\Pr(x|y)} . \end{aligned} \quad (5.2)$$

In the nomenclature of [20],  $\mathbf{ei}(X \rightarrow y)$  can be understood as the “total causal power” the system exerts when transitioning into state  $y$ .

2. The  $\phi$  measure aims to quantify a system’s irreducibility to disjoint parts, and the second step is to quantify how much of the total causal power isn’t accounted for by the disjoint parts (partition)  $\mathbf{P}$ . To do this, they define the *effective information beyond partition  $\mathbf{P}$* ,

$$\mathbf{ei}(X \rightarrow y/\mathbf{P}) \equiv D_{\text{KL}} \left[ \Pr(X|y) \parallel \prod_{i=1}^m \Pr(X_i^{\mathbf{P}} | y_i^{\mathbf{P}}) \right] . \quad (5.3)$$

The intuition behind  $\mathbf{ei}(X \rightarrow y/\mathbf{P})$  is that it quantifies the amount of causal power in  $\mathbf{ei}(X \rightarrow y)$  that is irreducible to the parts  $\mathbf{P}$  operating independently.<sup>2</sup>

3. After defining the causal power beyond an arbitrary partition  $\mathbf{P}$ , the third step is to find the partition that accounts for as much causal power as possible. This partition is called the *Minimum Information Partition*, or MIP. They define the MIP for a given state  $y$  as,<sup>3</sup>

$$\text{MIP}(y) \equiv \underset{\mathbf{P}}{\text{argmin}} \frac{\mathbf{ei}(X \rightarrow y/\mathbf{P})}{(m-1) \cdot \min_i H(X_i^{\mathbf{P}})} . \quad (5.4)$$

Finding the MIP of a system by brute force is incredibly computationally expensive, as it requires enumerating all partitions of  $n$  nodes scales  $O(n!)$  and even for supercomputers becomes intractable for  $n > 32$  nodes.

4. Fourth and finally, the system’s causal irreducibility when transitioning into state  $y \in Y$ ,  $\phi(y)$ , is the effective information beyond  $y$ ’s MIP,

$$\phi(y) \equiv \mathbf{ei}(X \rightarrow y/\mathbf{P} = \text{MIP}(y)) .$$

---

<sup>2</sup>In [3] they deviated slightly from this formulation, using a process termed “perturbing the wires”. However, subsequent work[32, 33] disavowed perturbing the wires and thus we don’t use it here. For discussion see Appendix 5.C

<sup>3</sup>In [3] they additionally consider the *total partition* as a special case. However, subsequent work[32, 33] disavowed the total partition and thus we don’t use it here.

### 5.3.1 Stateless $\phi$ is $\langle\phi\rangle$

In [3]  $\phi$  is defined for every state  $y \in Y$ , and a single system can have a wide range of  $\phi$ -values. In [5], they found this medley of state-dependent  $\phi$ -values unwieldy, and they decided to get a single number per system by averaging the effective information over all states  $y$ . This gives rise to the four corresponding stateless measures:

$$\begin{aligned}
\langle\mathbf{ei}(Y)\rangle &\equiv \mathbb{E}_y \mathbf{ei}(X \rightarrow y) = I(X:Y) \\
\langle\mathbf{ei}(Y/\mathbf{P})\rangle &\equiv \mathbb{E}_y \mathbf{ei}(X \rightarrow y/\mathbf{P}) = I(X:Y) - \sum_{i=1}^m I(X_i^{\mathbf{P}}:Y_i^{\mathbf{P}}) \\
\langle\text{MIP}\rangle &\equiv \underset{\mathbf{P}}{\operatorname{argmin}} \frac{\langle\mathbf{ei}(Y/\mathbf{P})\rangle}{(m-1) \cdot \min_i H(X_i^{\mathbf{P}})} \\
\langle\phi\rangle &\equiv \left\langle \mathbf{ei}(Y/\mathbf{P} = \langle\text{MIP}\rangle) \right\rangle.
\end{aligned} \tag{5.5}$$

Although the distinction has yet to affect qualitative results, researchers should note that  $\langle\phi\rangle \neq \mathbb{E}_Y \phi(y)$ . This is because whereas each  $y$  state can have a different MIP, for  $\langle\phi\rangle$  there's only one MIP for all states.

## 5.4 Room for improvement in $\phi$

$\phi(y)$  **can exceed**  $H(X)$ . Figure 5.1 shows examples OR-GET and OR-XOR. On average, each looks fine—they each have  $H(X) = 2$ ,  $I(X:Y) = 1.5$ , and  $\langle\phi\rangle = 1.189$  bits—nothing peculiar. This changes when examining the individual states  $y \in Y$ .

For OR-GET, the  $\phi(y = 10) \approx 2.58$  bits. Therefore  $\phi(y)$  *exceeds* the entropy of the entire system,  $H(XY) = H(X) = 2$  bits. This means that for  $y = 10$ , the “irreducible causal power” exceeds not just the total causal power,  $\mathbf{ei}(X \rightarrow y)$ , but  $\mathbf{ei}$ 's upperbound of  $H(X)$ ! This is concerning.

For OR-XOR,  $\phi(y = 11) \approx 1.08$  bits. This does not exceed  $H(X)$ , but it does exceed  $I(X:y = 11) = 1$  bit. Per eq. (5.5), in expectation  $\langle\mathbf{ei}(Y/\mathbf{P})\rangle \leq I(X:Y)$  for any partition  $\mathbf{P}$ . An information-theoretic interpretation of the state-dependent case would be more natural if likewise  $\mathbf{ei}(X \rightarrow y/\mathbf{P}) \leq I(X:Y = y)$  for any partition  $\mathbf{P}$ . Note this issue is not due simply to normalizing in eq. (5.4). For OR-GET and OR-XOR there's only one possible partition, and thus the normalization has no effect. The oddity arises from the equation for the effective information beyond a partition, eq. (5.3).

$\phi$  **sometimes decreases with duplicate computation.** In Figure 5.2 we take a simple system, AND-ZERO, and duplicate the AND gate yielding AND-AND. We see the two systems remain exceedingly similar. Each has  $H(X) = 2$  and  $I(X:Y) = 0.811$  bits. Likewise, each has two  $Y$  states occurring with probability  $3/4$  and  $1/4$ , giving  $\mathbf{ei}(X \rightarrow y)$  equal to 0.42 and 2.00 bits, respectively. However, their  $\phi$  values are quite different.

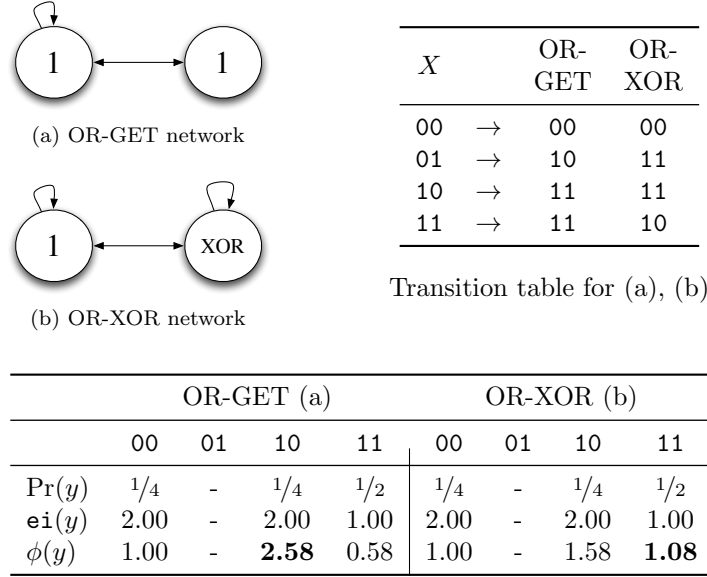


Figure 5.1: Example OR-GET shows that  $\phi(y)$  can exceed not only  $\mathbf{ei}(X \rightarrow y)$ , but  $H(X)$ ! A dash means that particular  $y$  is unreachable for the network. The concerning  $\phi$  values are **bolded**.

If we only knew that the  $\phi$ 's for AND-AND and AND-ZERO were different, we'd expect AND-AND's  $\phi$  to be higher because an AND node “does more” than a ZERO node (simply shutting off). But instead we get the opposite—AND-AND's highest  $\phi$  is *less* than AND-ZERO's lowest  $\phi$ ! An ideal measure of integrated information might be invariant or increase with duplicate computation, but it certainly wouldn't decrease!

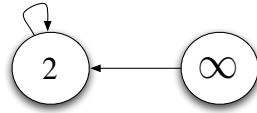
**$\phi$  does not increase with cooperation among diverse parts.** The  $\phi$  measure is often thought of as corresponding to the juxtaposition of “functional segregation” and “functional integration”. In a similar vein,  $\phi$  is also intuited as corresponding to “interdependence/cooperation among diverse parts”. Figure 5.3 presents four examples showing that these similar intuitions are not captured by the existing  $\phi$  measure.

In the first example, SHIFT (Figure 5.3a), each bit one step clockwise—nothing more, nothing less. The nodes are homogeneous and each node is determined by its preceding node.

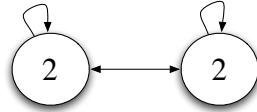
In the three remaining networks (Figures 5.3b–d), each node is a function of all nodes in its network (including itself). This is to maximize interdependence among the nodes, making the network highly “functionally integrated”. Having established high cooperation, we increase the diversity/“functional segregation” from Figure 5.3b to 5.3d.

By the aforementioned intuitions, we'd expect SHIFT (Figure 5.3a) to have the lowest  $\phi$  and 4321 (Figure 5.3d) to have the highest. But this is not the case. Instead, SHIFT, the network with the *least* cooperation (every node is a function of one other) and the *least* diverse mechanisms (all nodes have threshold 1) has a  $\phi$  far exceeding the others; SHIFT's lowest  $\phi$  value at 2.00 bits dwarfs the  $\phi$  values in Figures 5.3b–d.





(a) AND-ZERO network



(b) AND-AND network

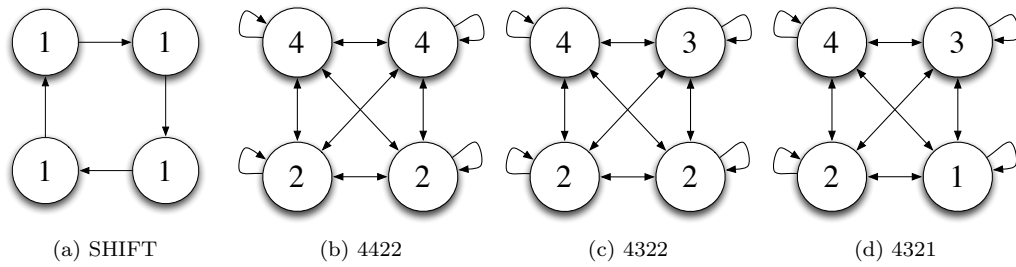
$X$		AND-ZERO	AND-AND
00	→	00	00
01	→	00	00
10	→	00	00
11	→	10	11

Transition table for (a), (b)

	AND-ZERO (a)				AND-AND (b)			
	00	01	10	11	00	01	10	11
$\text{Pr}(y)$	3/4	-	1/4	-	3/4	-	-	1/4
$\text{ei}(y)$	0.42	-	2.00	-	0.42	-	-	2.00
$\phi(y)$	0.33	-	1.00	-	0.25	-	-	0.00

Figure 5.2: Examples AND-ZERO and AND-AND show that  $\phi(y)$  sometimes *decreases* with duplicate computation. Here, the highest  $\phi$  of AND-AND is *less* than the lowest  $\phi$  of AND-ZERO. This carries into the average case with AND-ZERO's  $\langle \phi \rangle = 0.5$  and AND-AND's  $\langle \phi \rangle = 0.189$  bits. A dash means that particular  $y$  is unreachable for the network.

SHIFT having the highest  $\phi$  is unexpected, but it's not outright absurd. In SHIFT each node is wholly determined by an external force (the preceding node); so in some sense SHIFT is “integrated”. Whether it makes sense for SHIFT to have the highest integrated information ultimately comes down to precisely what is meant by the term “integration”. But even accepting that SHIFT is in some sense integrated, network 4321 is integrated for a stronger sense of the term. Therefore, until there's some argument that the awareness of SHIFT *should* be higher than 4321, from a purely theoretical perspective it makes sense to prefer 4321 over SHIFT.



Network	$I(X:Y)$	$\min_y \phi(y)$	$\max_y \phi(y)$	$\langle \phi \rangle$
SHIFT	4.000	2.000	2.000	2.000
4422	1.198	0.000	0.673	0.424
4322	1.805	0.322	1.586	1.367
4321	2.031	0.322	1.682	1.651

Figure 5.3: State-dependent  $\phi$  and  $\langle \phi \rangle$  tell the same story—the  $\phi$  value of SHIFT (a) trounces the  $\phi$  of the other three networks. Neither  $\phi$  measure is representative of cooperation among diverse parts.

## 5.5 A Novel Measure of Irreducibility to a Partition

Our proposed measure  $\psi$  quantifies the magnitude of information in  $I(X:y)$  (eq. (5.1)) that is irreducible to a partition of the system at time=0. We define our measure as,

$$\psi(\mathbf{X} : y) \equiv I(X:y) - \max_{\mathbf{P}} I_{\cup}(X_1^{\mathbf{P}}, \dots, X_m^{\mathbf{P}} : y) , \quad (5.6)$$

where  $\mathbf{P}$  enumerates over all partitions of set  $\mathbf{X}$ , and  $I_{\cup}$  is the information about state  $y$  conveyed by the “union” across the  $m$  parts at time=0. To compute the union information  $I_{\cup}$  we use the Partial Information Decomposition (PID) framework. In PID,  $I_{\cup}$  is the dual to  $I_{\cap}$ ; they are related by the inclusion–exclusion principle. Thus we can express  $I_{\cup}$  by,

$$I_{\cup}(X_1^{\mathbf{P}}, \dots, X_m^{\mathbf{P}} : y) = \sum_{\mathbf{S} \subseteq \{X_1^{\mathbf{P}}, \dots, X_m^{\mathbf{P}}\}} (-1)^{|\mathbf{S}|+1} I_{\cap}(S_1, \dots, S_{|\mathbf{S}|} : y) .$$

Conceptually, the intersection information  $I_{\cap}(S_1, \dots, S_{|\mathbf{S}|} : y)$  is the magnitude of information about state  $y$  that is conveyed “redundantly” by each  $S_i \in \mathbf{S}$ . Although there currently remains some debate[7, 14] about what is the best measure of  $I_{\cap}$ , there’s consensus that the intersection information  $n$  arbitrary random variables  $Z_1, \dots, Z_n$  convey about state  $Y = y$  must satisfy the following properties:

- (**GP**) Global Positivity:  $I_{\cap}(Z_1, \dots, Z_n : y) \geq 0$  with equality if  $\Pr(y) = 0$  or  $\Pr(y) = 1$ .
- (**M<sub>0</sub>**) Weak Monotonicity:  $I_{\cap}(Z_1, \dots, Z_n, W : y) \leq I_{\cap}(Z_1, \dots, Z_n : y)$  with equality if there exists  $Z_i \in \{Z_1, \dots, Z_n\}$  such that  $H(Z_i|W) = 0$ .
- (**S<sub>0</sub>**) Weak Symmetry:  $I_{\cap}(Z_1, \dots, Z_n : y)$  is invariant under reordering  $Z_1, \dots, Z_n$ .
- (**SR**) Self-Redundancy:  $I_{\cap}(Z_1 : y) = I(Z_1 : y) = D_{\text{KL}}[\Pr(Z_1|y) \parallel \Pr(Z_1)]$ . The intersection information a single predictor  $Z_1$  conveys about the target state  $Y = y$  is equal to the “specific surprise”[10] between the predictor and the target state.
- (**Eq**) Equivalence-Class Invariance:  $I_{\cap}(Z_1, \dots, Z_n : y)$  is invariant under substituting  $Z_i$  (for any  $i = 1, \dots, n$ ) by an informationally equivalent random variable<sup>4</sup>. [14] Similarly,  $I_{\cap}(Z_1, \dots, Z_n : y)$  is invariant under substituting state  $y$  for state  $w$  if  $\Pr(w|y) = \Pr(y|w) = 1$ .

Instead of choosing a particular  $I_{\cap}$  that satisfies the above properties, we will simply use these properties directly to bound the range of possible  $\psi$  values. Leveraging (**M<sub>0</sub>**), (**S<sub>0</sub>**), and (**SR**),

---

<sup>4</sup>Meaning that  $I_{\cap}$  is invariant under substituting  $Z_i$  with  $W$  if  $H(Z_i|W) = H(W|Z_i) = 0$ .

eq. (5.6) simplifies to,<sup>5</sup>

$$\begin{aligned}\psi(\mathbf{X} : y) &= I(X : y) - \max_{A \subset \mathbf{X}} I_{\cup}(A, B : y) \\ &= I(X : y) - \max_{A \subset \mathbf{X}} [I(A : y) + I(B : y) - I_{\cap}(A, B : y)] ,\end{aligned}\tag{5.7}$$

where  $A \neq \emptyset$  and  $B \equiv \mathbf{X} \setminus A$ .

From eq. (5.7), the only term left to be defined is  $I_{\cap}(A, B : y)$ . Leveraging **(GP)**, **(M<sub>0</sub>)**, and **(SR)**, we can bound this by  $0 \leq I_{\cap}(A, B : y) \leq \min [I(A : y), I(B : y)]$ .

Finally, we bound  $\psi$  by plugging in the above bounds on  $I_{\cap}(A, B : y)$  into eq. (5.7). With some algebra and leveraging assumption **(B)**, this becomes,<sup>6</sup>

$$\begin{aligned}\psi_{\min}(\mathbf{X} : y) &= \min_{A \subset \mathbf{X}} D_{\text{KL}} [\Pr(X_{1\dots n}|y) \parallel \Pr(A|y) \Pr(B|y)] \\ \psi_{\max}(\mathbf{X} : y) &= \min_{i \in \{1, \dots, n\}} D_{\text{KL}} [\Pr(X_{1\dots n}|y) \parallel \Pr(X_i) \Pr(X_{\sim i}|y)] ,\end{aligned}\tag{5.8}$$

where  $X_{\sim i}$  is the random variable of all nodes in  $X$  excluding node  $i$ . Then,  $\psi_{\min}(\mathbf{X} : y) \leq \psi(\mathbf{X} : y) \leq \psi_{\max}(\mathbf{X} : y)$ .

### 5.5.1 Stateless $\psi$ is $\langle \psi \rangle$

Matching how  $\langle \phi \rangle$  is defined in Section 5.3.1, to compute  $\langle \psi \rangle$  we weaken the properties in Section 5.5 so that they only apply to the average case, i.e., the properties **(GP)**, **(M<sub>0</sub>)**, **(S<sub>0</sub>)**, **(SR)**, and **(Eq)** don't have to apply for each  $I_{\cap}(Z_1, \dots, Z_n : y)$ , but merely for the average case  $I_{\cap}(Z_1, \dots, Z_n : Y)$ .

Via the same algebra<sup>7</sup>,  $\langle \psi \rangle$  simplifies to,

$$\begin{aligned}\langle \psi \rangle(X_1, \dots, X_n : Y) &\equiv I(X : Y) - \max_{\mathbf{P}} I_{\cup}(X_1^{\mathbf{P}}, \dots, X_m^{\mathbf{P}} : Y) \\ &= I(X : Y) - \max_{A \subset \mathbf{X}} I_{\cup}(A, B : Y) \\ &= I(X : Y) - \max_{A \subset \mathbf{X}} [I(A : Y) + I(B : Y) - I_{\cap}(A, B : Y)] ,\end{aligned}\tag{5.9}$$

where  $A \neq \emptyset$  and  $B \equiv \mathbf{X} \setminus A$ . Using the weakened properties, we have  $0 \leq I_{\cap}(A, B : Y) \leq \min [I(A : Y), I(B : Y)]$ . Plugging in these  $I_{\cap}$  bounds, we achieve the analogous bounds on  $\langle \psi \rangle$ ,<sup>8</sup>

$$\begin{aligned}\langle \psi \rangle_{\min}(\mathbf{X} : Y) &= \min_{A \subset \mathbf{X}} I(A : B | Y) \\ \langle \psi \rangle_{\max}(\mathbf{X} : Y) &= \min_{i \in \{1, \dots, n\}} D_{\text{KL}} [\Pr(X, Y) \parallel \Pr(X_{\sim i}, Y) \Pr(X_i)] ,\end{aligned}\tag{5.10}$$

<sup>5</sup>See Appendix 5.B.1 for a proof.

<sup>6</sup>See Appendix 5.B.2 for proofs.

<sup>7</sup>See Appendix 5.B.1 for a proof.

<sup>8</sup>See Appendix 5.B.3 for proofs.

where  $X_{\sim i}$  is the random variable of all nodes in  $X$  excluding node  $i$ . Then,  
 $\langle \psi \rangle_{\min}(\mathbf{X} : Y) \leq \langle \psi \rangle(\mathbf{X} : Y) \leq \langle \psi \rangle_{\max}(\mathbf{X} : Y)$ .

## 5.6 Contrasting $\psi$ versus $\phi$

**Theoretical benefits.** The overarching theoretical benefit of  $\psi$  is that it is entrenched within the rigorous Partial Information Decomposition framework[36]. PID builds a theoretically principled irreducibility measure from a redundancy measure  $I_{\cap}$ . Here we only take the most accepted properties of  $I_{\cap}$  to bound  $\psi$  from above and below. As the complexity community converges on the additional properties  $I_{\cap}$  must satisfy[7, 14], the derived bounds on  $\psi$  will contract.

The first benefit of  $\psi$ 's principled underpinning is that whereas  $\phi(y)$  can exceed the entropy of the whole system, i.e.,  $\phi(y) \not\leq H(X)$ ,  $\psi(y)$  is bounded by specific-surprise, i.e.,  $\psi(y) \leq I(X:y) = D_{\text{KL}} \left[ \Pr(X|y) \parallel \Pr(X) \right]$ . This gives  $\psi$  the natural info-theoretic interpretation for the state-dependent case which  $\phi$  lacks. A second benefit is that PID provides justification for  $\psi$  not needing a MIP normalization, and thus eliminates a longstanding concern about  $\phi$ [2]. The third benefit is that PID is a flexible framework that enables quantifying irreducibility to overlapping parts should we decide to explore it<sup>9</sup>.

One final perk is that  $\psi$  is substantially faster to compute. Whereas computing  $\phi$  scales<sup>10</sup>  $O(n!)$ , computing  $\psi$  scales<sup>11</sup>  $O(2^n)$ —a substantial improvement that may improve even further as the complexity community converges on additional properties of  $I_{\cap}$ .

**Practical differences.** The first row in Figure 5.4 shows two ways a network can be irreducible to atomic elements (the nodes) yet still reducible to disjoint parts. Compare AND-ZERO (Figure 5.4g) to AND-ZERO+KEEP (Figure 5.4a). Although AND-ZERO is irreducible, AND-ZERO+KEEP reduces to the bipartition separating the AND-ZERO component and the KEEP node. This reveals how fragile measures like  $\psi$  and  $\phi$  are—add a single disconnected node and they plummet to zero. Example 2x AND-ZERO (Figure 5.4b) shows that a reducible system can be composed entirely of irreducible parts.

Example KEEP-KEEP (Figure 5.4c) highlights the only known relative drawback of  $\psi$ : its current upperbound<sup>12</sup> is painfully loose. The desired irreducibility for KEEP-KEEP is zero bits, and indeed,  $\psi_{\min}$  is 0 bits, but  $\psi_{\max}$  is a monstrous 1 bit! We rightly expect tighter bounds for such easy examples like KEEP-KEEP. Tighter bounds on  $I_{\cap}$  (and thus  $\psi$ ) is an area of active research but as-is the bounds are loose.

Example GET-GET (Figure 5.4d) epitomizes the most striking difference between  $\psi$  and  $\phi$ .

<sup>9</sup>Technically there are multiple irreducibilities to overlapping parts as, unlike disjoint parts, the maximum union information over two overlapping parts is not equal to the maximum union information over  $m$  overlapping parts.

<sup>10</sup>This comes from eq. (5.4) enumerating all partitions (Bell's number) of  $n$  elements.

<sup>11</sup>This comes from eq. (5.7) enumerating all  $2^{n-1} - 1$  bipartitions of  $n$  elements.

<sup>12</sup>The current upperbounds are  $\psi_{\max}$  in eq. (5.8) and  $\langle \psi \rangle_{\max}$  in eq. (5.10).

By property **(Eq)**, the  $\psi$  values for KEEP-KEEP and GET-GET are provably equal (making the desired  $\psi$  for GET-GET also zero bits), yet their  $\phi$  values couldn't be more different. Although  $\phi$  agrees KEEP-KEEP is zero,  $\phi$  firmly places GET-GET at the maximal (!) two bits of irreducibility. Whereas  $\psi$  views GET nodes akin to a system-wide KEEP,  $\phi$  views GET nodes as highly integrative.

The primary benefit of making KEEPs and GETs equivalent is that  $\psi$  is zero for chains of GETs such as the SHIFT network (Fig. 5.3a). This enables  $\psi$  to better match our intuition for “cooperation among diverse parts”. For example, in Figure 5.3 the network with the highest  $\phi$  is the counter-intuitive SHIFT; on the other hand, the network with the highest  $\psi$  is the more sensible 4321 (see bottom table in Figure 5.4).

The third row in Figure 5.4 shows a difference related to KEEPs vs GETs—how  $\psi$  and  $\phi$  respectively treat self-connections. In ANDtriplet (Figure 5.4e) each node integrates information about two nodes. Likewise, in iso-ANDtriplet (Figure 5.4f) each node integrates information about two nodes, but the information is about *itself* and one other.

Just as  $\psi$  views KEEP and GET nodes equivalently,  $\psi$  views self and cross connections equivalently. In fact, by property **(Eq)** the  $\psi$  values for ANDtriplet and iso-ANDtriplet are provably equal. On the other hand,  $\phi$  considers self and cross connections differently in that  $\phi$  can only decrease when adding a self-connection. As such, the  $\phi$  for iso-ANDtriplet is less than ANDtriplet.

The fourth row in Figure 5.4 shows this same self-connections business carrying over to duplicate computations. Although AND-AND (Figure 5.4h) and AND-ZERO (Figure 5.4g) perform the same computation, AND-AND has an additional self-connection that pushes AND-AND's  $\phi$  below that of AND-ZERO. By **(Eq)**,  $\psi$  is provably invariant under such duplicate computations.

## 5.7 Conclusion

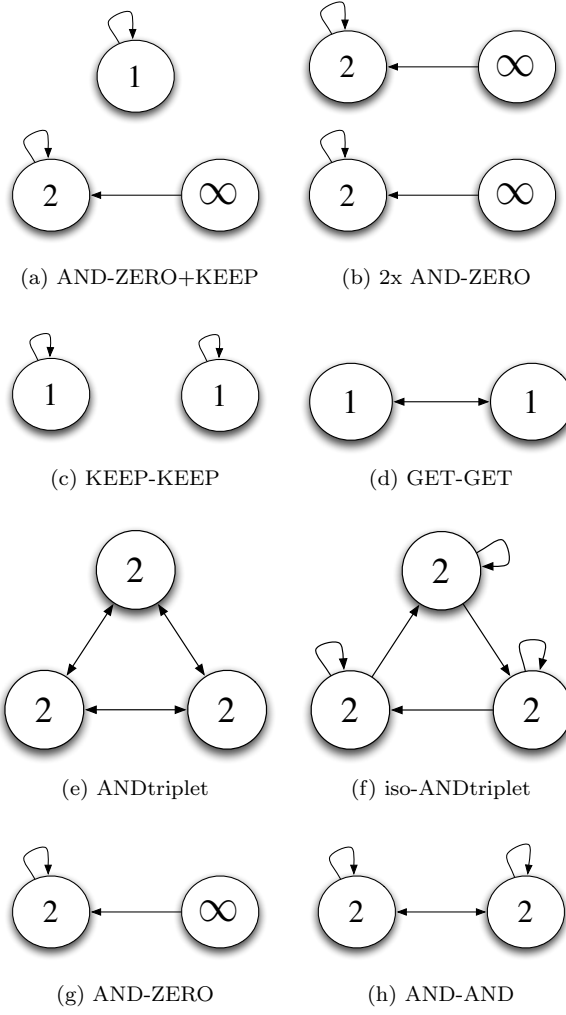
Regardless of any connection to consciousness, and purely as a measure of functional irreducibility, we have three concerns about  $\phi$ : (1) state-dependent  $\phi$  can exceed the entropy of the entire system; (2)  $\phi$  often decreases with duplicate computation; (3)  $\phi$  doesn't match the intuition of “cooperation among diverse parts”.

We introduced a new irreducibility measure,  $\psi$ , that solves all three concerns and otherwise stays close to the original spirit of  $\phi$ —i.e., the quantification of a system's irreducibility to disjoint parts. Based in Partial Information Decomposition,  $\psi$  has other desirable properties, such as not needing a MIP normalization and being substantially faster to compute.

Finally, we contrasted  $\psi$  versus  $\phi$  with simple, concrete examples.

Although we recommend using  $\psi$  over  $\phi$ , the  $\psi$  measure remains imperfect. The most notable areas for improvement are:

1. The current  $\psi$  bounds are too loose. We need to tighten the  $I_\cap$  bounds, which will tighten the



Network	$I(X:Y)$	$\langle\phi\rangle$	$\langle\psi\rangle_{\min}$	$\langle\psi\rangle_{\max}$
AND-ZERO+KEEP (a)	1.81	0	0	0.50
2x AND-ZERO (b)	1.62	0	0	0.50
KEEP-KEEP (c)	2.00	0	0	1.00
GET-GET (d)	2.00	2.00	0	1.00
SHIFT (Fig. 5.3a)	4.00	2.00	0	1.00
4422 (Fig. 5.3b)	1.20	0.42	0.33	0.50
4322 (Fig. 5.3c)	1.81	1.37	0.68	0.88
4321 (Fig. 5.3d)	2.03	1.65	0.78	1.00
ANDtriplet (e)	2.00	2.00	0.16	0.75
iso-ANDtriplet (f)	2.00	1.07	0.16	0.75
AND-ZERO (g)	0.81	0.50	0.19	0.5
AND-AND (h)	0.81	0.19	0.19	0.5

Figure 5.4: Contrasting  $\langle\phi\rangle$  versus  $\langle\psi\rangle$  for exemplary networks.

derived bounds on  $\psi$  and  $\langle\psi\rangle$ .

2. Justify why a measure of conscious experience should privilege irreducibility to disjoint parts over irreducibility to overlapping parts.
3. Reformalize the work on qualia in [4] using  $\psi$ .
4. Although not specific to  $\psi$ , there needs to be a stronger justification for the chosen distribution on  $X$ .

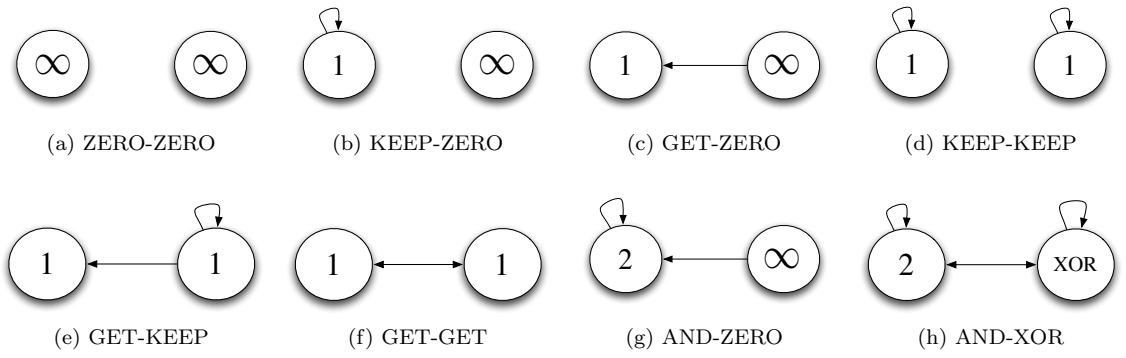


# Appendix

## 5.A Reading the network diagrams

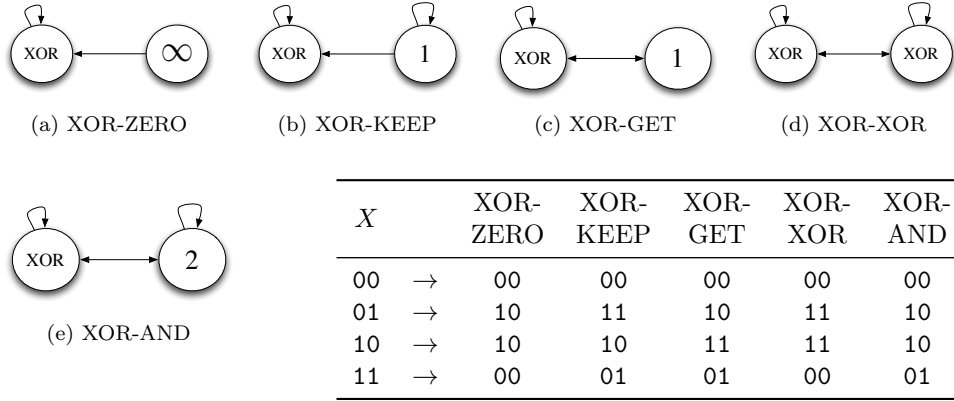
We present eight doublet networks and their transition tables so you can see how the network diagram specifies the transition table. Figure 5.5 shows eight network diagrams to build your intuition. The number inside each node is that node's *activation threshold*. A node updates to 1 (conceptually an “ON”) if there are at least as many of inputs ON as its activation threshold; e.g. a node with an inscribed 2 updates to a 1 if two or more incoming wires are ON. An activation threshold of  $\infty$  means the node always updates to 0 (conceptually an “OFF”). A binary string denotes the state of the network, read left to right.

We take the AND-ZERO network (Figure 5.5g) as an example. Although the AND-ZERO network can never output 01 or 11 (Fig. 1b), we still consider states 01, 11 as equally possible states at time=0. This is because  $X_0$  is uniformly distributed per assumption **(A)**. The state of the AND-node (left) at time=1 is a function of *both* nodes at time=0. For example, in the AND-ZERO gate, the left binary digit is the state of the AND-node and the right binary digit is the state of the ZERO-node.



$X$		ZERO-ZERO	KEEP-ZERO	GET-ZERO	KEEP-KEEP	GET-KEEP	GET-GET	AND-ZERO	AND-XOR
00	→	00	00	00	00	00	00	00	00
01	→	00	00	10	01	11	10	00	01
10	→	00	10	00	10	00	01	00	01
11	→	00	10	10	11	11	11	10	10

Figure 5.5: Eight doublet networks with transition tables.



Network	$I(X:Y)$	$\langle\phi\rangle$	$\langle\psi\rangle_{\min}$	$\langle\psi\rangle_{\max}$
ZERO-ZERO (Fig. 5.5a)	0	0	0	0
KEEP-ZERO (Fig. 5.5b)	1.0	0	0	0
KEEP-KEEP (Fig. 5.5d)	2.0	0	0	1.0
GET-ZERO (Fig. 5.5c)	1.0	1.0	0	0
GET-KEEP (Fig. 5.5e)	1.0	0	0	0
GET-GET (Fig. 5.5f)	2.0	2.0	0	1.0
AND-ZERO (Fig. 5.2a)	0.811	0.5	0.189	0.5
AND-KEEP	1.5	0.189	0	0.5
AND-GET	1.5	1.189	0	0.5
AND-AND (Fig. 5.2b)	0.811	0.189	0.189	0.5
AND-XOR (Fig. 5.5h)	1.5	1.189	0.5	1.0
XOR-ZERO (a)	1.0	1.0	1.0	1.0
XOR-KEEP (b)	2.0	1.0	0	1.0
XOR-GET (c)	2.0	2.0	0	1.0
XOR-AND (e)	1.5	1.189	0.5	1.0
XOR-XOR (d)	1.0	1.0	1.0	1.0

Figure 5.6: Networks, transition tables, and measures for the diagnostic doublets.

## 5.B Necessary proofs

### 5.B.1 Proof that the max union of bipartitions covers all partitions

**Lemma 5.B.1.** *Given properties  $(\mathbf{S}_0)$  and  $(\mathbf{M}_0)$ , the maximum union information conveyed by a partition of predictors  $\mathbf{X} = \{X_1, \dots, X_n\}$  about state  $Y = y$  equals the maximum union information conveyed by a bipartition of  $\mathbf{X}$  about state  $Y = y$ .*

*Proof.* We prove that the maximum Information conveyed by a Partition,  $\text{lbDp}(\mathbf{X} : y)$ , equals Information conveyed by a Bipartition,  $\text{lbB}(\mathbf{X} : y)$  by showing,

$$\text{lbDp}(\mathbf{X} : y) \leq \text{lbB}(\mathbf{X} : y) \leq \text{lbDp}(\mathbf{X} : y) . \quad (5.11)$$

We first show that  $\text{lbB}(\mathbf{X} : y) \leq \text{lbDp}(\mathbf{X} : y)$ . By their definitions,

$$\begin{aligned} \text{lbDp}(\mathbf{X} : y) &\equiv \max_{\mathbf{P}} I_{\cup}(\mathbf{P} : y) \\ \text{lbB}(\mathbf{X} : y) &\equiv \max_{|\mathbf{P}|=2} I_{\cup}(\mathbf{P} : y) , \end{aligned} \quad (5.12)$$

where  $\mathbf{P}$  enumerates over all partitions of set  $\mathbf{X}$ .

By removing the restriction that  $|\mathbf{P}| = 2$  from the maximization in  $\text{lbB}$  we arrive at  $\text{lbDp}$ . As removing a restriction can only increase the maximum, thus  $\text{lbB}(\mathbf{X} : y) \leq \text{lbDp}(\mathbf{X} : y)$ .

We next show that  $\text{lbDp}(\mathbf{X} : y) \leq \text{lbB}(\mathbf{X} : y)$ , meaning we must show that,

$$\max_{\mathbf{P}} I_{\cup}(\mathbf{P} : y) \leq \max_{|\mathbf{P}|=2} I_{\cup}(\mathbf{P} : y) . \quad (5.13)$$

Without loss of generality, we choose an arbitrary subset/part  $S \subset \mathbf{X}$ . This yields the bipartition of parts  $\{S, \mathbf{X} \setminus S\}$ . We then further partition the second part,  $\mathbf{X} \setminus S$ , into  $k$  disjoint subcomponents denoted  $\{T_1, \dots, T_k\}$  where  $2 \leq k \leq n - |S|$ , creating an arbitrary partition  $\mathbf{P} = \{S, T_1, \dots, T_k\}$ . We now need to show that,

$$I_{\cup}(S, T_1, \dots, T_k : y) \leq I_{\cup}(S, \mathbf{X} \setminus S : y) .$$

By  $(\mathbf{M}_0)$  equality condition, we can append each subcomponent  $T_1, \dots, T_k$  to  $\{S, \mathbf{X} \setminus S\}$  without changing the union-information, because every subcomponent  $T_i \preceq \mathbf{X} \setminus S$ . Then applying  $(\mathbf{S}_0)$ , we re-order the parts so that  $S, T_1, \dots, T_k$  come first. This yields,

$$I_{\cup}(S, T_1, \dots, T_k : y) \leq I_{\cup}(S, T_1, \dots, T_k, \mathbf{X} \setminus S : y) .$$

Applying  $(\mathbf{M}_0)$  inequality condition, adding the predictor  $\mathbf{X} \setminus S$  can only increase the union information. Therefore we prove eq. (5.13), which proves eq. (5.11), that  $\text{lbDp}(\mathbf{X} : y) = \text{lbB}(\mathbf{X} : y)$ .  $\square$

### 5.B.2 Bounds on $\psi(X_1, \dots, X_n : y)$

**Lemma 5.B.2.** *Given  $(\mathbf{M}_0)$ ,  $(\mathbf{SR})$  and the predictors  $X_1, \dots, X_n$  are independent, i.e.  $H(X_{1\dots n}) = \sum_{i=1}^n H(X_i)$ , then,*

$$\psi(X_1, \dots, X_n : y) \leq \min_{i \in \{1, \dots, n\}} D_{\text{KL}} \left[ \Pr(X_{1\dots n} | y) \parallel \Pr(X_i) \Pr(X_{1\dots n \setminus i} | y) \right] .$$

*Proof.* Applying  $(\mathbf{M}_0)$  inequality condition, we have  $I_{\cap}(A, B : y) \leq \min [I(A : y), I(B : y)]$ . Via the inclusion-exclusion rule, this entails  $I_{\cup}(A, B : y) \geq \max [I(A : y), I(B : y)]$ , and we use this to upper-bound  $\psi(X_1, \dots, X_n : y)$ . The random variable  $A \neq \emptyset$ ,  $B \equiv \mathbf{X} \setminus A$ , and  $AB \equiv X_{1\dots n}$ .

$$\begin{aligned} \psi(X_1, \dots, X_n : y) &= I(X_{1\dots n} : y) - \max_{A \subset \mathbf{X}} I_{\cup}(A, B : y) \\ &\leq I(X_{1\dots n} : y) - \max_{A \subset \mathbf{X}} \max [I(A : y), I(B : y)] \end{aligned}$$

By symmetry of complementary bipartitions, every  $B$  will be an  $A$  at some point. So we can drop the  $B$  term.

$$= I(X_{1\dots n} : y) - \max_{A \subset \mathbf{X}} I(A : y) .$$

For two random variables  $A$  and  $A'$  such that  $A \preceq A'$ ,  $I(A : y) \leq I(A' : y)$ .<sup>13</sup> Therefore, there will always be a maximizing subset of  $\mathbf{X}$  with size  $n - 1$ .

$$\begin{aligned} \psi(X_1, \dots, X_n : y) &\leq I(X_{1\dots n} : y) - \max_{\substack{A \subset \mathbf{X} \\ |A|=n-1}} I(A : y) \\ &= I(X_{1\dots n} : y) - \max_{i \in \{1, \dots, n\}} I(X_{1\dots n \setminus i} : y) \\ &= \min_{i \in \{1, \dots, n\}} I(X_{1\dots n} : y) - I(X_{1\dots n \setminus i} : y) \\ &= \min_{i \in \{1, \dots, n\}} I(X_i : y | X_{1\dots n \setminus i}) \\ &= \min_{i \in \{1, \dots, n\}} D_{\text{KL}} \left[ \Pr(X_{1\dots n} | y) \parallel \Pr(X_i | X_{1\dots n \setminus i}) \Pr(X_{1\dots n \setminus i} | y) \right] . \end{aligned}$$

Now applying that the predictors  $\mathbf{X}$  are independent,  $\Pr(x_i | x_{1\dots n \setminus i}) = \Pr(x_i)$ . This yields,

$$\psi(X_1, \dots, X_n : y) \leq \min_{i \in \{1, \dots, n\}} D_{\text{KL}} \left[ \Pr(X_{1\dots n} | y) \parallel \Pr(X_i) \Pr(X_{1\dots n \setminus i} | y) \right] .$$

---

<sup>13</sup> $I(A : y) \leq I(A' : y)$  because  $I(A' : y) = I(A : y) + I(A' : y | A)$ .



**Lemma 5.B.3.** *Given (GP), (SR) and predictors  $X_1, \dots, X_n$  are independent, i.e.  $H(X_{1..n}) = \sum_{i=1}^n H(X_i)$ , then,*

$$\begin{aligned} \psi(X_1, \dots, X_n : y) &\geq \min_{A \subset \mathbf{X}} I(A : B | y) \\ &= \min_{A \subset \mathbf{X}} D_{\text{KL}} \left[ \Pr(X_{1..n} | y) \parallel \Pr(A | y) \Pr(B | y) \right]. \end{aligned}$$

*Proof.* First, from the definition of  $I_{\cup}$ ,  $I_{\cup}(A, B : y) = I(A : y) + I(B : y) - I_{\cap}(A, B : y)$ . Then applying (GP), we have  $I_{\cup}(A, B : y) \leq I(A : y) + I(B : y)$ . We use this to lowerbound  $\psi(X_1, \dots, X_n : y)$ . The random variable  $A \neq \emptyset$ ,  $B \equiv \mathbf{X} \setminus A$ , and  $AB \equiv X_{1..n}$ .

$$\begin{aligned} \psi(X_1, \dots, X_n : y) &= I(X_{1..n} : y) - \max_{A \subset \mathbf{X}} I_{\cup}(A, B : y) \\ &\geq I(X_{1..n} : y) - \max_{A \subset \mathbf{X}} [I(A : y) + I(B : y)] \\ &= \min_{A \subset \mathbf{X}} I(AB : y) - I(A : y) - I(B : y) \\ &= \min_{A \subset \mathbf{X}} I(A : y | B) - I(A : y) \\ &= \min_{A \subset \mathbf{X}} D_{\text{KL}} \left[ \Pr(AB | y) \parallel \Pr(B | y) \Pr(A | B) \right] - D_{\text{KL}} [\Pr(A | y) \parallel \Pr(A)] \\ &= \min_{A \subset \mathbf{X}} \sum_{a,b} \Pr(ab | y) \log \frac{\Pr(ab | y)}{\Pr(b | y) \Pr(a | b)} + \sum_a \Pr(a | y) \log \frac{\Pr(a)}{\Pr(a | y)}. \end{aligned}$$

We now add  $\sum_b \Pr(b | ay)$  in front of the right-most  $\sum_a$ . We can do this because  $\sum_b \Pr(b | ay) = 1.0$ . This then yields,

$$\begin{aligned} \psi(X_1, \dots, X_n : y) &\geq \min_{A \subset \mathbf{X}} \sum_{a,b} \Pr(ab | y) \log \frac{\Pr(ab | y)}{\Pr(b | y) \Pr(a | b)} + \Pr(b | ay) \Pr(a | y) \log \frac{\Pr(a)}{\Pr(a | y)} \\ &= \min_{A \subset \mathbf{X}} \sum_{a,b} \Pr(ab | y) \left[ \log \frac{\Pr(ab | y)}{\Pr(b | y) \Pr(a | b)} + \log \frac{\Pr(a)}{\Pr(a | y)} \right] \\ &= \min_{A \subset \mathbf{X}} \sum_{a,b} \Pr(ab | y) \log \frac{\Pr(ab | y) \Pr(a)}{\Pr(a | y) \Pr(b | y) \Pr(a | b)}. \end{aligned}$$

Now applying that the predictors  $\mathbf{X}$  are independent,  $\Pr(a | b) = \Pr(a)$ ; thus we can cancel  $\Pr(a)$



for  $\Pr(a|b)$ . This yields,

$$\begin{aligned}\psi(X_1, \dots, X_n : y) &\geq \min_{A \subset \mathbf{X}} \sum_{a,b} \Pr(ab|y) \log \frac{\Pr(ab|y)}{\Pr(a|y) \Pr(b|y)} \\ &= \min_{A \subset \mathbf{X}} D_{\text{KL}} \left[ \Pr(X_{1\dots n}|y) \parallel \Pr(A|y) \Pr(B|y) \right] .\end{aligned}$$

□

### 5.B.3 Bounds on $\langle \psi \rangle(X_1, \dots, X_n : Y)$

**Lemma 5.B.4.** *Given  $(\mathbf{M}_0)$ ,  $(\mathbf{SR})$  and the predictors  $X_1, \dots, X_n$  are independent, i.e.  $H(X_{1\dots n}) = \sum_{i=1}^n H(X_i)$ , then,*

$$\langle \psi \rangle(X_1, \dots, X_n : Y) \leq \min_{i \in \{1, \dots, n\}} D_{\text{KL}} \left[ \Pr(X_{1\dots n}, Y) \left\| \Pr(X_{1\dots n \setminus i}, Y) \Pr(X_i) \right. \right] .$$

*Proof.* First, using the same reasoning in Lemma 5.B.2, we have,

$$\begin{aligned} \langle \psi \rangle(\mathbf{X} : Y) &\leq I(X_{1\dots n} : Y) - \max_{i \in \{1, \dots, n\}} I(X_{1\dots n \setminus i} : Y) \\ &= \min_{i \in \{1, \dots, n\}} I(X_{1\dots n} : Y) - I(X_{1\dots n \setminus i} : Y) \\ &= \min_{i \in \{1, \dots, n\}} I(X_i : Y | X_{1\dots n \setminus i}) \\ &= \min_{i \in \{1, \dots, n\}} D_{\text{KL}} \left[ \Pr(X_{1\dots n}, Y) \left\| \Pr(X_i | X_{1\dots n \setminus i}) \Pr(X_{1\dots n \setminus i}, Y) \right. \right] . \end{aligned}$$

Now applying that the predictors  $\mathbf{X}$  are independent,  $\Pr(X_i | X_{1\dots n \setminus i}) = \Pr(X_i)$ . This yields,

$$\langle \psi \rangle(\mathbf{X} : Y) \leq \min_{i \in \{1, \dots, n\}} D_{\text{KL}} \left[ \Pr(X_{1\dots n}, Y) \left\| \Pr(X_{1\dots n \setminus i}, Y) \Pr(X_i) \right. \right] .$$

□

**Lemma 5.B.5.** *Given  $(\mathbf{GP})$ ,  $(\mathbf{SR})$  and predictors  $X_1, \dots, X_n$  are independent, i.e.  $H(X_{1\dots n}) = \sum_{i=1}^n H(X_i)$ , then,*

$$\langle \psi \rangle(X_1, \dots, X_n : Y) \geq \min_{A \subset \mathbf{X}} I(A : B | Y) .$$

*Proof.* First, using the same reasoning in Lemma 5.B.3, we have,

$$\begin{aligned} \langle \psi \rangle(X_1, \dots, X_n : Y) &\geq I(X_{1\dots n} : Y) - \max_{A \subset \mathbf{X}} [I(A : Y) + I(B : Y)] \\ &= \min_{A \subset \mathbf{X}} I(AB : Y) - I(A : Y) - I(B : Y) \\ &= \min_{A \subset \mathbf{X}} I(A : B | Y) - I(A : B) . \end{aligned}$$

Now applying that the predictors  $\mathbf{X}$  are independent,  $I(A:B) = 0$ . This yields,

$$\langle \psi \rangle(X_1, \dots, X_n : Y) \geq \min_{A \subset \mathbf{X}} I(A:B|Y) \ .$$

□

## 5.C Definition of intrinsic $\text{ei}(y/\mathbf{P})$ a.k.a. “perturbing the wires”

State-dependent  $\text{ei}$  across a partition, fully written as  $\text{ei}(X \rightarrow y/\mathbf{P})$  and abbreviated  $\text{ei}(y/\mathbf{P})$ , is defined by eq. (5.14). The probability distribution of the “intrinsic information” in the entire system,  $\Pr(X \rightarrow y)$ , is simply  $\Pr(X|y)$  (eq. (5.15)).<sup>14</sup>

$$\text{ei}(X \rightarrow y/\mathbf{P}) \equiv D_{\text{KL}} \left[ \Pr^*(X \rightarrow y) \left\| \prod_{i=1}^m \Pr(X_i^{\mathbf{P}} \rightarrow y_i^{\mathbf{P}}) \right\| \right] \quad (5.14)$$

$$= D_{\text{KL}} \left[ \Pr(X|y) \left\| \prod_{i=1}^m \Pr^*(X_i^{\mathbf{P}}|y_i^{\mathbf{P}}) \right\| \right]. \quad (5.15)$$

Balduzzi/Tononi [3] define the probability distribution describing the intrinsic information from the whole system  $X$  to state  $y$  as,

$$\Pr(X \rightarrow y) = \Pr(X|y) = \left\{ \Pr(x|y) : \forall x \in X \right\}.$$

They then define probability distribution describing the intrinsic information from a part  $X_i^{\mathbf{P}}$  to a state  $y_i^{\mathbf{P}}$  as,

$$\Pr^*(X_i^{\mathbf{P}} \rightarrow y_i^{\mathbf{P}}) \equiv \Pr^*(X_i^{\mathbf{P}}|Y_i^{\mathbf{P}} = y_i^{\mathbf{P}}) = \left\{ \Pr^*(x_i^{\mathbf{P}}|y_i^{\mathbf{P}}) : \forall x_i^{\mathbf{P}} \in X_i^{\mathbf{P}} \right\}.$$

First we define the fundamental property of the  $\Pr^*$  distribution. Given a state  $x_i^{\mathbf{P}}$ , the probability of a state  $y_i^{\mathbf{P}}$  is computed by probability that each node in the state  $y_i^{\mathbf{P}}$  independently reaches the state specified by  $y_i^{\mathbf{P}}$ ,

$$\Pr^*(y_i^{\mathbf{P}}|x_i^{\mathbf{P}}) \equiv \prod_{j=1}^{|\mathbf{P}_i|} \Pr(y_{i,j}^{\mathbf{P}}|x_i^{\mathbf{P}}). \quad (5.16)$$

Then we define the join distribution relative to eq. (5.16):

$$\begin{aligned} \Pr^*(x_i^{\mathbf{P}}, y_i^{\mathbf{P}}) &= \Pr^*(x_i^{\mathbf{P}}) \Pr^*(y_i^{\mathbf{P}}|x_i^{\mathbf{P}}) \\ &= \Pr^*(x_i^{\mathbf{P}}) \prod_{j=1}^{|\mathbf{P}_i|} \Pr(y_{i,j}^{\mathbf{P}}|x_i^{\mathbf{P}}). \end{aligned}$$

Then applying assumption **(B)**,  $X$  follows a discrete uniform distribution, so  $\Pr^*(x_i^{\mathbf{P}}) \equiv \Pr(x_i^{\mathbf{P}}) =$

---

<sup>14</sup>It's worth nothing that  $\Pr^*(X|y) \neq \Pr(X|y)$ .

$1/|X_i^{\mathbf{P}}|$ . This gives us the complete definition of  $\Pr^*(x_i^{\mathbf{P}}, y_i^{\mathbf{P}})$ ,

$$\Pr^*(x_i^{\mathbf{P}}, y_i^{\mathbf{P}}) = \Pr(x_i^{\mathbf{P}}) \prod_{j=1}^{|\mathbf{P}_i|} \Pr(y_{i,j}^{\mathbf{P}} | x_i^{\mathbf{P}}) . \quad (5.17)$$

With the joint  $\Pr^*$  distribution defined, we can compute anything we want by summing over the eq. (5.17)—such as the expressions for  $\Pr^*(y_i^{\mathbf{P}})$  and  $\Pr^*(x_i^{\mathbf{P}} | y_i^{\mathbf{P}})$ ,

$$\begin{aligned} \Pr^*(y_i^{\mathbf{P}}) &= \sum_{x_i^{\mathbf{P}} \in X_i^{\mathbf{P}}} \Pr^*(x_i^{\mathbf{P}}, y_i^{\mathbf{P}}) \\ \Pr^*(x_i^{\mathbf{P}} | y_i^{\mathbf{P}}) &= \frac{\Pr^*(x_i^{\mathbf{P}}, y_i^{\mathbf{P}})}{\Pr^*(y_i^{\mathbf{P}})} . \end{aligned}$$

## 5.D Misc proofs

Given the properties **(GP)**, **(SR)**, and the predictors  $X_1, \dots, X_n$  are independent, i.e.  $H(X_{1\dots n}) = \sum_{i=1}^n H(X_i)$ , we show that  $\langle \phi \rangle \not\leq \langle \psi \rangle$ . This is equivalent to,

$$\langle \psi \rangle_{\min}(\mathbf{X} : Y) \leq \langle \phi \rangle(\mathbf{X} : Y) .$$

*Proof.* We prove the above by showing that for any bipartition  $\mathbf{P}$ , that  $\langle \psi \rangle_{\min}(\mathbf{X} : Y) \leq \langle \mathbf{ei}(Y/\mathbf{P}) \rangle$ .

For a bipartition  $\mathbf{P}$ ,

$$\begin{aligned} \langle \psi \rangle_{\min}(\mathbf{X} : Y) &= I(X_1^{\mathbf{P}} : X_2^{\mathbf{P}} | Y) \\ &= I(X_1^{\mathbf{P}} : X_2^{\mathbf{P}} | Y) - I(X_1^{\mathbf{P}} : X_2^{\mathbf{P}}) \\ &= I(X : Y) - I(X_1^{\mathbf{P}} : Y) - I(X_2^{\mathbf{P}} : Y) \\ \langle \mathbf{ei}(Y/\mathbf{P}) \rangle &= I(X : Y) - I(X_1^{\mathbf{P}} : Y_1^{\mathbf{P}}) - I(X_1^{\mathbf{P}} : Y_1^{\mathbf{P}}) . \end{aligned}$$

$$\begin{aligned} \langle \mathbf{ei}(Y/\mathbf{P}) \rangle - \langle \psi \rangle_{\min}(\mathbf{X} : Y) &= I(X : Y) - I(X_1^{\mathbf{P}} : Y_1^{\mathbf{P}}) - I(X_1^{\mathbf{P}} : Y_1^{\mathbf{P}}) - I(X : Y) + I(X_1^{\mathbf{P}} : Y) + I(X_2^{\mathbf{P}} : Y) \\ &= I(X_1^{\mathbf{P}} : Y) - I(X_1^{\mathbf{P}} : Y_1^{\mathbf{P}}) + I(X_2^{\mathbf{P}} : Y) - I(X_2^{\mathbf{P}} : Y_2^{\mathbf{P}}) \\ &= I(X_1^{\mathbf{P}} : Y_2^{\mathbf{P}} | Y_1^{\mathbf{P}}) + I(X_2^{\mathbf{P}} : Y_1^{\mathbf{P}} | Y_2^{\mathbf{P}}) \\ &\geq 0 . \end{aligned}$$

And we complete the proof that  $\langle \psi \rangle_{\min} \leq \langle \phi \rangle$ . Therefore,  $\langle \phi \rangle \not\leq \langle \psi \rangle$ .

□

## 5.E Setting $t = 1$ without loss of generality

Given  $t$  stationary surjective functions that may be different or the same, denoted  $f_1 \cdots f_t$ , we define the state of system at time  $t$ , denoted  $X_t$ , as the application of the  $t$  functions to the state of the system at time 0, denoted  $X_0$ ,

$$X_t = f_t \left( f_{t-1} \left( \cdots f_2 (f_1 (X_0)) \cdots \right) \right) .$$

We instantiate an empty “dictionary function”  $g(\bullet)$ . Then for every  $x_0 \in X_0$  we assign,

$$g(x_0) \equiv f_t \left( f_{t-1} \left( \cdots f_2 (f_1 (x_0)) \cdots \right) \right) . \quad \forall x_0 \in X_0$$

At the end of this process we have a function  $g$  that accomplishes any chain of stationary functions  $f_1 \cdots f_t$  in a single step for the entire domain of  $f_1$ . So instead of studying the transformation,

$$X_0 \xrightarrow{f_1 \cdots f_t} X_t ,$$

we can equivalently study the transformation,

$$X_0 \xrightarrow{g} Y .$$

Here’s an example using mechanism  $f_1 = f_2 = f_3 = f_4 = \text{AND-GET}$ .

time=0		$t = 1$		$t = 2$		$t = 3$		$t = 4$
00	→	00	→	00	→	00	→	00
01	→	00	→	00	→	00	→	00
10	→	01	→	00	→	00	→	00
11	→	11	→	11	→	10	→	00
$g(\bullet)$		AND-GET		AND-AND		AND-ZERO		ZERO-ZERO

Table 5.1: Applying the update rule “AND-GET”, over four timesteps.

# Bibliography

- [1] Dimitris Anastassiou. Computational analysis of the synergy among multiple interacting genes. *Molecular Systems Biology*, 3:83, 2007.
- [2] David Balduzzi. personal communication.
- [3] David Balduzzi and Giulio Tononi. Integrated information in discrete dynamical systems: motivation and theoretical framework. *PLoS Computational Biology*, 4(6):e1000091, Jun 2008.
- [4] David Balduzzi and Giulio Tononi. Qualia: The geometry of integrated information. *PLoS Computational Biology*, 5(8), 2009.
- [5] Adam B. Barrett and Anil K. Seth. Practical measures of integrated information for time-series data. *PLoS Computational Biology*, 2010.
- [6] Anthony J. Bell. The co-information lattice. In S. Amari, A. Cichocki, S. Makino, and N. Murata, editors, *Fifth International Workshop on Independent Component Analysis and Blind Signal Separation*. Springer, 2003.
- [7] Nils Bertschinger, Johannes Rauh, Eckehard Olbrich, and Jürgen Jost. Shared information – new insights and problems in decomposing information in complex systems. *CoRR*, abs/1210.5902, 2012.
- [8] N. J. Cerf and C. Adami. Negative entropy and information in quantum mechanics. *Phys. Rev. Lett.*, 79:5194–5197, Dec 1997.
- [9] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. John Wiley, New York, NY, 1991.
- [10] M. R. DeWeese and M. Meister. How to measure the information gained from one symbol. *Network*, 10:325–340, Nov 1999.
- [11] T. G. Dietterich, S. Becker, and Z. Ghahramani, editors. *Group Redundancy Measures Reveal Redundancy Reduction in the Auditory Pathway*, Cambridge, MA, 2002. MIT Press.



- [12] Itay Gat and Naftali Tishby. Synergy and redundancy among brain cells of behaving monkeys. In *Advances in Neural Information Proceedings systems*, pages 465–471. MIT Press, 1999.
- [13] Timothy J. Gawne and Barry J. Richmond. How independent are the messages carried by adjacent inferior temporal cortical neurons? *Journal of Neuroscience*, 13:2758–71, 1993.
- [14] V. Griffith, E. K. P. Chong, R. G. James, C. J. Ellison, and J. P. Crutchfield. Intersection information based on common randomness. *ArXiv e-prints*, October 2013.
- [15] Virgil Griffith and Christof Koch. Quantifying synergistic mutual information. In M. Prokopenko, editor, *Guided Self-Organization: Inception*. Springer, 2014.
- [16] Peter Gács and Jack Körner. Common information is far less than mutual information. *Problems of Control and Information Theory*, 2(2):149–162, 1973.
- [17] Te Sun Han. Nonnegative entropy measures of multivariate symmetric correlations. *Information and Control*, 36(2):133–156, 1978.
- [18] Malte Harder, Christoph Salge, and Daniel Polani. A bivariate measure of redundant information. *CoRR*, abs/1207.2080, 2012.
- [19] A. Jakulin and I. Bratko. Analyzing attribute dependencies. In *Lecture Notes in Artificial Intelligence*, pages 229–240, 2838 2003.
- [20] Kevin B Korb, Lucas R Hope, and Erik P Nyberg. Information-theoretic causal power. In *Information Theory and Statistical Learning*, pages 231–265. Springer, 2009.
- [21] Peter E. Latham and Sheila Nirenberg. Synergy, redundancy, and independence in population codes, revisited. *Journal of Neuroscience*, 25(21):5195–5206, May 2005.
- [22] Hua Li and Edwin K. P. Chong. On a connection between information and group lattices. *Entropy*, 13(3):683–708, 2011.
- [23] Joseph T. Lizier, Benjamin Flecker, and Paul L. Williams. Towards a synergy-based approach to measuring information modification. *CoRR*, abs/1303.3440, 2013.
- [24] W. J. McGill. Multivariate information transmission. *Psychometrika*, 19:97–116, 1954.
- [25] S Nirenberg, S M Carcieri, A L Jacobs, and P E Latham. Retinal ganglion cells act largely as independent encoders. *Nature*, 411(6838):698–701, Jun 2001.
- [26] Sheila Nirenberg and Peter E. Latham. Decoding neuronal spike trains: How important are correlations? *Proceedings of the National Academy of Sciences*, 100(12):7348–7353, 2003.

- [27] S Panzeri, A Treves, S Schultz, and E T Rolls. On decoding the responses of a population of neurons from short time windows. *Neural Comput*, 11(7):1553–1577, Oct 1999.
- [28] G Pola, A Thiele, K P Hoffmann, and S Panzeri. An exact method to quantify the information transmitted by different mechanisms of correlational coding. *Network*, 14(1):35–60, Feb 2003.
- [29] E. Schneidman, W. Bialek, and M.J. Berry II. Synergy, redundancy, and independence in population codes. *Journal of Neuroscience*, 23(37):11539–53, 2003.
- [30] Elad Schneidman, Susanne Still, Michael J. Berry, and William Bialek. Network information and connected correlations. *Phys. Rev. Lett.*, 91(23):238701–238705, Dec 2003.
- [31] Giulio Tononi. An information integration theory of consciousness. *BMC Neuroscience*, 5(42), November 2004.
- [32] Giulio Tononi. Consciousness as integrated information: a provisional manifesto. *Biological Bulletin*, 215(3):216–242, Dec 2008.
- [33] Giulio Tononi. The integrated information theory of consciousness: An updated account. *Archives Italiennes de Biologie*, 150(2/3):290–326, 2012.
- [34] Vinay Varadan, David M. Miller, and Dimitris Anastassiou. Computational inference of the molecular logic for synaptic connectivity in *c. elegans*. *Bioinformatics*, 22(14):e497–e506, 2006.
- [35] Eric W. Weisstein. Antichain. <http://mathworld.wolfram.com/Antichain.html>, 2011.
- [36] Paul L. Williams and Randall D. Beer. Nonnegative decomposition of multivariate information. *CoRR*, abs/1004.2515, 2010.
- [37] Stefan Wolf and Jürg Wullschleger. Zero-error information and applications in cryptography. *Proc IEEE Information Theory Workshop*, 04:1–6, 2004.
- [38] A. D. Wyner. The common information of two dependent random variables. *IEEE Transactions in Information Theory*, 21(2):163–179, March 1975.